

SWABEY OGILVY RENAULT
RECEIVED

APR 14 1999

PATENT COOPERATION TREATY

WO 99/15639
PCT/CA98/00884

M. 8 | 9 | 10 | 11 | 12 | 1 | 2 | 3 | 4 | 5 | 6 PCT
P.M.

**NOTICE INFORMING THE APPLICANT OF THE
COMMUNICATION OF THE INTERNATIONAL
APPLICATION TO THE DESIGNATED OFFICES**

(PCT Rule 47.1(c), first sentence)

From the INTERNATIONAL BUREAU

To:

COTE, France
Swabey Ogilvy Renault
Suite 1600
1981 McGill College Avenue
Montréal, Québec H3A 2Y3
CANADA

Date of mailing (day/month/year) 01 April 1999 (01.04.99)		IMPORTANT NOTICE	
Applicant's or agent's file reference 1770-182PCT			
International application No. PCT/CA98/00884	International filing date (day/month/year) 18 September 1998 (18.09.98)	Priority date (day/month/year) 19 September 1997 (19.09.97)	
Applicant MCGILL UNIVERSITY et al			

1. Notice is hereby given that the International Bureau has communicated, as provided in Article 20, the international application to the following designated Offices on the date indicated above as the date of mailing of this Notice:
AU,BR,CN,EP,IL,JP,KP,KR,US

In accordance with Rule 47.1(c), third sentence, those Offices will accept the present Notice as conclusive evidence that the communication of the international application has duly taken place on the date of mailing indicated above and no copy of the international application is required to be furnished by the applicant to the designated Office(s).

2. The following designated Offices have waived the requirement for such a communication at this time:
AL,AM,AP,AT,AZ,BA,BB,BG,BY,CA,CH,CU,CZ,DE,DK,EA,EE,ES,FI,GB,GE,GH,GM,HR,HU,ID,IS,KE,KG,KZ,LC,LK,LR,LS,LT,LU,LV,MD,MG,MK,MN,MW,MX,NO,NZ,OA,PL,PT,RO,RU,SD,SE,SG,SI,SK,SL,TJ,TM,TR,TT,UA,UG,UZ,VN,YU,ZW
The communication will be made to those Offices only upon their request. Furthermore, those Offices do not require the applicant to furnish a copy of the international application (Rule 49.1(a-bis)).

3. Enclosed with this Notice is a copy of the international application as published by the International Bureau on 01 April 1999 (01.04.99) under No. WO 99/15639

REMINDER REGARDING CHAPTER II (Article 31(2)(a) and Rule 54.2)

If the applicant wishes to postpone entry into the national phase until 30 months (or later in some Offices) from the priority date, a demand for international preliminary examination must be filed with the competent International Preliminary Examining Authority before the expiration of 19 months from the priority date.

It is the applicant's sole responsibility to monitor the 19-month time limit.

Note that only an applicant who is a national or resident of a PCT Contracting State which is bound by Chapter II has the right to file a demand for international preliminary examination.

REMINDER REGARDING ENTRY INTO THE NATIONAL PHASE (Article 22 or 39(1))

If the applicant wishes to proceed with the international application in the national phase, he must, within 20 months or 30 months, or later in some Offices, perform the acts referred to therein before each designated or elected Office.

For further important information on the time limits and acts to be performed for entering the national phase, see the Annex to Form PCT/IB/301 (Notification of Receipt of Record Copy) and Volume II of the PCT Applicant's Guide.

The International Bureau of WIPO 34, chemin des Colombettes 1211 Geneva 20, Switzerland	Authorized officer J. Zahra
Facsimile No. (41-22) 740.14.35	Telephone No. (41-22) 338.83.38

Continuation of Form PCT/IB/308

NOTICE INFORMING THE APPLICANT OF THE COMMUNICATION OF
THE INTERNATIONAL APPLICATION TO THE DESIGNATED OFFICES

Date of mailing (day/month/year) 01 April 1999 (01.04.99)	IMPORTANT NOTICE
Applicant's or agent's file reference 1770-182PCT	International application No. PCT/CA98/00884
<p>The applicant is hereby notified that, at the time of establishment of this Notice, the time limit under Rule 46.1 for making amendments under Article 19 has not yet expired and the International Bureau had received neither such amendments nor a declaration that the applicant does not wish to make amendments.</p>	

PATENT COOPERATION TREATY

PCT

NOTIFICATION OF ELECTION

(PCT Rule 61.2)

From the INTERNATIONAL BUREAU

To:

United States Patent and Trademark
Office
(Box PCT)
Crystal Plaza 2
Washington, DC 20231
ÉTATS-UNIS D'AMÉRIQUE

in its capacity as elected Office

Date of mailing (day/month/year)

28 May 1999 (28.05.99)

International application No.

PCT/CA98/00884

Applicant's or agent's file reference

1770-182PCT

International filing date (day/month/year)

18 September 1998 (18.09.98)

Priority date (day/month/year)

19 September 1997 (19.09.97)

Applicant

ROULEAU, Guy, A. et al

1. The designated Office is hereby notified of its election made:



in the demand filed with the International Preliminary Examining Authority on:

15 April 1999 (15.04.99)



in a notice effecting later election filed with the International Bureau on:

2. The election ☒ was

was not

made before the expiration of 19 months from the priority date or, where Rule 32 applies, within the time limit under Rule 32.2(b).

The International Bureau of WIPO
34, chemin des Colombettes
1211 Geneva 20, Switzerland

Facsimile No.: (41-22) 740.14.35

Authorized officer

C. Carrié

Telephone No.: (41-22) 338.83.38

PATENT COOPERATION TREATY

PCT

INTERNATIONAL PRELIMINARY EXAMINATION REPORT

(PCT Article 36 and Rule 70)

REC'D 29 DEC 1999

WIPO PCT

Applicant's or agent's file reference 1770-182PCT		FOR FURTHER ACTION See Notification of Transmittal of International Preliminary Examination Report (Form PCT/IPEA/416)
International application No. PCT/CA98/00884	International filing date (day/month/year) 18/09/1998	Priority date (day/month/year) 19/09/1997
International Patent Classification (IPC) or national classification and IPC C12N15/00		
Applicant McGILL UNIVERSITY et al.		



1. This international preliminary examination report has been prepared by this International Preliminary Examining Authority and is transmitted to the applicant according to Article 36.
2. This REPORT consists of a total of 6 sheets, including this cover sheet.

☒ This report is also accompanied by ANNEXES, i.e. sheets of the description, claims and/or drawings which have been amended and are the basis for this report and/or sheets containing rectifications made before this Authority (see Rule 70.16 and Section 607 of the Administrative Instructions under the PCT).

 These annexes consist of a total of 3 sheets.

3. This report contains indications relating to the following items:

- I ☒ Basis of the report
- II ☐ Priority
- III ☐ Non-establishment of opinion with regard to novelty, inventive step and industrial applicability
- IV ☐ Lack of unity of invention
- V ☒ Reasoned statement under Article 35(2) with regard to novelty, inventive step or industrial applicability; citations and explanations supporting such statement
- VI ☒ Certain documents cited
- VII ☐ Certain defects in the international application
- VIII ☒ Certain observations on the international application

Date of submission of the demand 15/04/1999	Date of completion of this report 17. 12. 99
Name and mailing address of the international preliminary examining authority:  European Patent Office D-80298 Munich Tel. +49 89 2399 - 0 Tx: 523656 epmu d Fax: +49 89 2399 - 4465	Authorized officer Merlos-Lange. A.M. Telephone No. +49 89 2399 8559 

**INTERNATIONAL PRELIMINARY
EXAMINATION REPORT**

International application No. PCT/CA98/00884

I. Basis of the report

1. This report has been drawn on the basis of *(substitute sheets which have been furnished to the receiving Office in response to an invitation under Article 14 are referred to in this report as "originally filed" and are not annexed to the report since they do not contain amendments.)*:

Description, pages:

1-24 as originally filed

Claims, No.:

1-12 as received on 09/10/1999 with letter of 05/10/1999

Drawings, sheets:

1/9-9/9 as originally filed

2. The amendments have resulted in the cancellation of:

- ☐ the description, pages:
☐ the claims, Nos.:
☐ the drawings, sheets:

3. ☐ This report has been established as if (some of) the amendments had not been made, since they have been considered to go beyond the disclosure as filed (Rule 70.2(c)):

4. Additional observations, if necessary:

see separate sheet

**INTERNATIONAL PRELIMINARY
EXAMINATION REPORT**

International application No. PCT/CA98/00884

V. Reasoned statement under Article 35(2) with regard to novelty, inventive step or industrial applicability; citations and explanations supporting such statement

1. Statement

Novelty (N)	Yes: Claims 1-12
	No: Claims
Inventive step (IS)	Yes: Claims 1-12
	No: Claims
Industrial applicability (IA)	Yes: Claims 1-6, 8-11
	No: Claims

2. Citations and explanations

see separate sheet

VI. Certain documents cited

1. Certain published documents (Rule 70.10)

and / or

2. Non-written disclosures (Rule 70.9)

see separate sheet

VIII. Certain observations on the international application

The following observations on the clarity of the claims, description, and drawings or on the question whether the claims are fully supported by the description, are made:

see separate sheet

**INTERNATIONAL PRELIMINARY
EXAMINATION REPORT - SEPARATE SHEET**

International application No. PCT/CA98/00884

- 1). Section I, Point 4, add. observation

The sequence listing is shown in additional sheets 1/4 to 4/4.

- 2). Section VI

The IPER is established on the opinion that the present application enjoys a valid priority. In case of an unvalid priority of 09.09.1997, document "Europ. J. of Human Genetics, January 1998, 6, 89-94, Philibert, R. A. et. al." may become relevant for the assessment of novelty of claim 1 when the application enters the regional phase.

- 3). For the assessment of the present claims 7 and 12 on the question whether they are industrially applicable, no unified criteria exist in the PCT Contracting States. The patentability can also be dependent upon the formulation of the claims. The EPO, for example, does not recognize as industrially applicable the subject-matter of claims to the use of a compound in medical treatment, but may allow, however, claims to a known compound for first use in medical treatment and the use of such a compound for the manufacture of a medicament for a new medical treatment.

Claims 7 and 12 relate to subject-matter considered by this Authority to be covered by the provisions of Rule 67.1(iv) PCT. Consequently, no opinion will be formulated with respect to the industrial applicability of the subject-matter of these claims (Article 34(4)(a)(i) PCT).

- 4). New claim 9 now refers to "a method of categorizing psychiatric patients according to their genotype to maximize response to treatment patients ...", which does however not appear to be supported in the original disclosure. With respect to the "use" claims 10 to 12 it is noted that they also appear broader than originally filed claims 5, 6 and 7 insofar as they are not dependent on these claims. Therefore the new claims are not limited to the use of determined allelic variants being obtained from a nucleic acid sample of a patient according to claim 4 (original claim 5) or to the use of a non-human mammal model for screening of therapeutic agents according to claim 7 (original claim 8).

In view of this, said claims are not considered to conform with the requirements of Art. 34 (2)(b) PCT.

5). Section VIII

When considering claim 1, it is not clear whether it is directed to a (wild type?) hGT1 gene comprising the sequence as set forth in Figs. 3 and 4A to 4C and containing transcribed polymorphic CAG repeat or whether it is directed to the particular allelic CAG repeat variants thereof. Furthermore, in the absence of the complete definition of the allelic CAG repeat variant as given on page 5, lines 26-32 or lines 11 to 16, the claim is not rendered more clear. The definition of "alleles -3, -2, -1, 0, and 1" is an arbitrary one introduced by the applicant and therefore meaningless to the skilled person unless the full meaning is included in the claims. Finally, it would appear that claim 1 refers to human GT1 (hGT1). However, reference to Fig. 3 which shows a human and a mouse GT1 sequence, introduces some doubt whether the mouse sequence should be involved in the scope of claim 1 or not.

In view of the above, the dependent claims 2-12 are not clearly defined in the sense of Art. 6 PCT as well.

With respect to claim 8 it is further noted that it appears to be incomplete (A method to identify genes **part of or interacting with** a biochemical ...). Moreover, the claimed method is not defined by particular procedure steps to identify a gene which forms part of or which interacts with a biochemical pathway affected by the hGT1 gene. Screening of samples with probes or primers derived from the (wild type?) hGT1 sequence does not appear to result in the identification of the desired gene which interacts for example with the biochemical pathway but rather in the identification of allelic variants of hGT1 gene. As already mentioned above, it is not clear whether claim 9 refers to the (wildtype) hGT1 gene of claim 1 or to particular allelic CAG repeat variants thereof.

6). Section V

None of the available prior art discloses or suggests means and methods as described in the present application which therefore appears to conform with the

**INTERNATIONAL PRELIMINARY
EXAMINATION REPORT - SEPARATE SHEET**

International application No. PCT/CA98/00884

requirements of novelty and inventive step according to Art. 33(2), (3) PCT.

(51) International Patent Classification ⁶ : C12N 15/00, C12Q 1/68, C07K 14/47, A01K 67/027	A1	(11) International Publication Number: WO 99/15639
		(43) International Publication Date: 1 April 1999 (01.04.99)

(81) **Designated States:** AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GE, GH, GM, HR, HU, ID, IL, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

(74) **Agents:** COTE, France et al.; Swabey Ogilvy Renault, Suite 1600, 1981 McGill College Avenue, Montréal, Québec H3A 2Y3 (CA).

*With international search report.
Before the expiration of the time limit for amending the
claims and to be republished in the event of the receipt of
amendments.*

(54) Title: POLYMORPHIC CAG REPEAT-CONTAINING GENE AND USES THEREOF

		130	140	150	160	170	180
D29801M (1>600)	→	AGGGCAGCCACTTTCCCAGCATTTCCGACCTTCCTTCCCTACCTCTCCCACTATGCCCAA					
GCT10D04 (1>320)		TCCTTCCCACCTCTCCACCTACTCTCTCGT					
		AGGGCAGCCACTTTCCCAGCATTTCCGACCTTCCTTCCCTACCTCTCCCACTATGCCAA					
		190	200	210	220	230	240
D29801M (1>600)	→	CAGTG-CAGGGTGTGGCGAGGGGGGCCACTCTCAAGAAGGTGCACGACCACTCTCG					
GCT10D04 (1>320)	↑	CTGTCCAGGGTGTGGCGAGGGGGCCNACTCTCTAAGAAGTGCACGACCACTCTCG					
	↑	CWGTGCAGGGTGTGGCGAGGGGGGCCACTCTCAAGAAGTGCACGACCACTCTCG					
		250	260	270	280	290	300
D29801M (1>600)	→	CAGCCTCATGATAGCGGCATGACTGCCAATGCGACCTTGGCTCAGGGGACAGGGTTCAG					
GCT10D04 (1>320)	↑	CAGCCCCATGACGAGCGGTGACTGCCAATGCGACCTTGGCTCAGGGGACAGGGTTCAG					
Oligo SCZ-15 (1>24)	↑	GGGGACAGGGGTTCAG					
	↑	CAGCCYCATGAYAGGCCGMGTASTGCCARYKCSARCCTGGCYCGRGGCACGGGTTCAG					

A

310 320 330 340 350 360

D29801M (1>600) → AATCTTTCAGCGTTACGAGCCTTGCCGCCCTTGCTACG
GCT1D0D4 (1>320) ← AATCTTTCAGCTACCAGTCGGCGCCCTCAGCTATGACCAGCAGCAGCAGCAGCAGCAG
Oligo SCZ-15 (1>24) → AATCTTC
AATCTTTCAYGCCTACCAAGTCGGCGCCCTTGCTATGACCAGcagcagcagcagcagcagc

370 380 390 400 410 420

D29801M (1>600) → GCAGCAGCAGCAGAAGCATTTACAGGCCCTCACACGAGGAAGAACTCTCAC
GCT1D0D4 (1>320) ← CAGCAGCAGCAGCAGCAGCAAGCCCTTCAGAGCGCCAGCATGCCAGGAAAACCTCCAT
CagcagcagcAgCAGCAGCAGCAAGCMCTTCARRGGCGRACCAYRCCAGGAACCHCEAY

430 440 450 460 470 480

D29801M (1>600) → TACCAGAAGCTCGCCAAGTACCAAACAATTGAGCAGAACGGCCAGGGCTACTGTCCA-CC
GCT1D0D4 (1>320) TACCAAAAGCTCGCCAAGTACGACGATACGGGGCAGCAGGCGAGGGCTACTG-CCAGCG
Oligo SCZ-16 (1>23) ← AGCATACGGGGCAGCAGGCGGAC

1

The present invention relates to hGT1 gene, a polymorphic CAG repeat-containing gene and its uses thereof for the diagnosis, prognosis and treatment of psychiatric diseases, such as schizophrenia.

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

INTERNATIONAL SEARCH REPORT

International Application No

PCT/CA 98/00884

A. CLASSIFICATION OF SUBJECT MATTER

IPC 6 C12N15/00 C12Q1/68 C07K14/47 A01K67/027

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 6 C07K C12Q

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	IMAI, Y. ET AL.: "Cloning of a retinoic acid-induced gene, GT1, in the embryonal carcinoma cell line P19: neuron-specific expression in the mouse brain" MOLECULAR BRAIN RESEARCH, XP002093328 see the whole document ---	1
A	ROBITAILLE, Y. ET AL.: "The neuropathology of CAG repeat diseases: review and update of genetic and molecular features" BRAIN PATHOLOGY, vol. 7, no. 3, July 1997, pages 901-926, XP002093329 see the whole document --- -/--	2

☒ Further documents are listed in the continuation of box C.

☒ Patent family members are listed in annex.

* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

"&" document member of the same patent family

Date of the actual completion of the international search

12 February 1999

Date of mailing of the international search report

25/02/1999

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Chambonnet, F

INTERNATIONAL SEARCH REPORT

International Application No
PCT/CA 98/00884

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	WO 95 01437 A (UNIV MINNESOTA) 12 January 1995 see the whole document ----	2-10
A	WO 97 18825 A (UNIV BRITISH COLUMBIA ;KALCHMAN MICHAEL (CA); HAYDEN MICHAEL R (CA) 29 May 1997 see the whole document ----	2-10
A	MACIEL, P. ET AL.: "Correlation between CAG repeat length and clinical features in Machado-Joseph Disease" AMERICAN JOURNAL OF HUMAN GENETICS, vol. 57, no. 1, July 1995, pages 54-61, XP002093330 ----	1,2,4,10
A	JOOPER, R. ET AL.: "Apolipoprotein E genotype in Schizophrenia" AMERICAN JOURNAL OF MEDICAL GENETICS (NEUROPSYCHIATRIC GENETICS), vol. 67, no. 2, 9 April 1996, page 235 XP002093331 see the whole document ----	2
P,X	PHILIBERT, R.A. ET AL.: "The characterization and sequence analysis of thirty CTG-repeat containing genomic cosmid clones" EUROPEAN JOURNAL OF HUMAN GENETICS, vol. 6, no. 1, January 1998, pages 89-94, XP002093332 see the whole document ----	1
T	TURECKI, G. ET AL.: "Schizophrenia and chromosome 6p" AMERICAN JOURNAL OF MEDICAL GENETICS, vol. 74, no. 2, 1997, pages 195-198, XP002093333 see the whole document -----	2

INTERNATIONAL SEARCH REPORT

Information on patent family members

International Application No

PCT/CA 98/00884

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
WO 9501437 A	12-01-1995	US 5834183 A	10-11-1998
		CA 2166117 A	12-01-1995
		EP 0707647 A	24-04-1996
		JP 9501049 T	04-02-1998
		US 5741645 A	21-04-1998
WO 9718825 A	29-05-1997	CA 2238075 A	29-05-1997
		EP 0873132 A	28-10-1998

RECEIVED

From the
INTERNATIONAL PRELIMINARY EXAMINING AUTHORITY

DEC 30 1999

To:

COTE, France
SWABEY OGILVY RENAULT
1981 McGill College Avenue
Suite 1600
Montréal, Québec H3A 2Y3
CANADAA.M. P.M.
PCT 9 10 11 12 1 2 3 4 5 6 7 8NOTIFICATION OF TRANSMITTAL OF
THE INTERNATIONAL PRELIMINARY
EXAMINATION REPORT
(PCT Rule 71.1)Date of mailing
(day/month/year)

17. 12. 99

Applicant's or agent's file reference
1770-182PCT

IMPORTANT NOTIFICATION

International application No.
PCT/CA98/00884International filing date (day/month/year)
18/09/1998Priority date (day/month/year)
19/09/1997Applicant
MCGILL UNIVERSITY et al.

1. The applicant is hereby notified that this International Preliminary Examining Authority transmits herewith the international preliminary examination report and its annexes, if any, established on the international application.
2. A copy of the report and its annexes, if any, is being transmitted to the International Bureau for communication to all the elected Offices.
3. Where required by any of the elected Offices, the International Bureau will prepare an English translation of the report (but not of any annexes) and will transmit such translation to those Offices.

4. REMINDER

The applicant must enter the national phase before each elected Office by performing certain acts (filing translations and paying national fees) within 30 months from the priority date (or later in some Offices) (Article 39(1)) (see also the reminder sent by the International Bureau with Form PCT/IB/301).

Where a translation of the international application must be furnished to an elected Office, that translation must contain a translation of any annexes to the international preliminary examination report. It is the applicant's responsibility to prepare and furnish such translation directly to each elected Office concerned.

For further details on the applicable time limits and requirements of the elected Offices, see Volume II of the PCT Applicant's Guide.

Name and mailing address of the IPEA/

European Patent Office
D-80298 Munich
Tel. +49 89 2399 - 0 Tx: 523656 epmu d
Fax: +49 89 2399 - 4465

Authorized officer

Vullo, C

Tel. +49 89 2399-8061





PATENT COOPERATION TREATY

PCT

INTERNATIONAL PRELIMINARY EXAMINATION REPORT

(PCT Article 36 and Rule 70)

Applicant's or agent's file reference 1770-182PCT		FOR FURTHER ACTION See Notification of Transmittal of International Preliminary Examination Report (Form PCT/IPEA/416)	
International application No. PCT/CA98/00884	International filing date (day/month/year) 18/09/1998	Priority date (day/month/year) 19/09/1997	
International Patent Classification (IPC) or national classification and IPC C12N15/00			
Applicant McGILL UNIVERSITY et al.			
<p>1. This international preliminary examination report has been prepared by this International Preliminary Examining Authority and is transmitted to the applicant according to Article 36.</p> <p>2. This REPORT consists of a total of 6 sheets, including this cover sheet.</p> <p><input checked="" type="checkbox"/> This report is also accompanied by ANNEXES, i.e. sheets of the description, claims and/or drawings which have been amended and are the basis for this report and/or sheets containing rectifications made before this Authority (see Rule 70.16 and Section 607 of the Administrative Instructions under the PCT).</p> <p>These annexes consist of a total of 3 sheets.</p>			
<p>3. This report contains indications relating to the following items:</p> <ul style="list-style-type: none"> I <input checked="" type="checkbox"/> Basis of the report II <input type="checkbox"/> Priority III <input type="checkbox"/> Non-establishment of opinion with regard to novelty, inventive step and industrial applicability IV <input type="checkbox"/> Lack of unity of invention V <input checked="" type="checkbox"/> Reasoned statement under Article 35(2) with regard to novelty, inventive step or industrial applicability; citations and explanations supporting such statement VI <input checked="" type="checkbox"/> Certain documents cited VII <input type="checkbox"/> Certain defects in the international application VIII <input checked="" type="checkbox"/> Certain observations on the international application 			
Date of submission of the demand 15/04/1999		Date of completion of this report 17. 12. 99	
Name and mailing address of the international preliminary examining authority:  European Patent Office D-80298 Munich Tel. +49 89 2399 10 Tx: 523656 epmu d Fax: +49 89 2399 - 4465		Authorized officer Merlos-Lange, A.M. Telephone No. +49 89 2399 8559 	

**INTERNATIONAL PRELIMINARY
EXAMINATION REPORT**

International application No. PCT/CA98/00884

I. Basis of the report

1. This report has been drawn on the basis of *(substitute sheets which have been furnished to the receiving Office in response to an invitation under Article 14 are referred to in this report as "originally filed" and are not annexed to the report since they do not contain amendments.)*:

Description, pages:

1-24 as originally filed

Claims, No.:

1-12 as received on 09/10/1999 with letter of 05/10/1999

Drawings, sheets:

1/9-9/9 as originally filed

2. The amendments have resulted in the cancellation of:

- ☐ the description, pages:
☐ the claims, Nos.:
☐ the drawings, sheets:

3. ☐ This report has been established as if (some of) the amendments had not been made, since they have been considered to go beyond the disclosure as filed (Rule 70.2(c)):

4. Additional observations, if necessary:

see separate sheet

**INTERNATIONAL PRELIMINARY
EXAMINATION REPORT**

International application No. PCT/CA98/00884

V. Reasoned statement under Article 35(2) with regard to novelty, inventive step or industrial applicability; citations and explanations supporting such statement

1. Statement

Novelty (N)	Yes:	Claims	1-12
	No:	Claims	
Inventive step (IS)	Yes:	Claims	1-12
	No:	Claims	
Industrial applicability (IA)	Yes:	Claims	1-6, 8-11
	No:	Claims	

2. Citations and explanations

see separate sheet

VI. Certain documents cited

1. Certain published documents (Rule 70.10)

and / or

2. Non-written disclosures (Rule 70.9)

see separate sheet

VIII. Certain observations on the international application

The following observations on the clarity of the claims, description, and drawings or on the question whether the claims are fully supported by the description, are made:

see separate sheet

**INTERNATIONAL PRELIMINARY
EXAMINATION REPORT - SEPARATE SHEET**

International application No. PCT/CA98/00884

- 1). Section I, Point 4, add. observation

The sequence listing is shown in additional sheets 1/4 to 4/4.

- 2). Section VI

The IPER is established on the opinion that the present application enjoys a valid priority. In case of an invalid priority of 09.09.1997, document "Europ. J. of Human Genetics, January 1998, 6, 89-94, Philibert, R. A. et. al." may become relevant for the assessment of novelty of claim 1 when the application enters the regional phase.

- 3). For the assessment of the present claims 7 and 12 on the question whether they are industrially applicable, no unified criteria exist in the PCT Contracting States. The patentability can also be dependent upon the formulation of the claims. The EPO, for example, does not recognize as industrially applicable the subject-matter of claims to the use of a compound in medical treatment, but may allow, however, claims to a known compound for first use in medical treatment and the use of such a compound for the manufacture of a medicament for a new medical treatment.

Claims 7 and 12 relate to subject-matter considered by this Authority to be covered by the provisions of Rule 67.1(iv) PCT. Consequently, no opinion will be formulated with respect to the industrial applicability of the subject-matter of these claims (Article 34(4)(a)(i) PCT).

- 4). New claim 9 now refers to "a method of categorizing psychiatric patients according to their genotype to maximize response to treatment patients ...", which does however not appear to be supported in the original disclosure. With respect to the "use" claims 10 to 12 it is noted that they also appear broader than originally filed claims 5, 6 and 7 insofar as they are not dependent on these claims. Therefore the new claims are not limited to the use of determined allelic variants being obtained from a nucleic acid sample of a patient according to claim 4 (original claim 5) or to the use of a non-human mammal model for screening of therapeutic agents according to claim 7 (original claim 8).

In view of this, said claims are not considered to conform with the requirements of Art. 34 (2)(b) PCT.

5). Section VIII

When considering claim 1, it is not clear whether it is directed to a (wild type?) hGT1 gene comprising the sequence as set forth in Figs. 3 and 4A to 4C and containing transcribed polymorphic CAG repeat or whether it is directed to the particular allelic CAG repeat variants thereof. Furthermore, in the absence of the complete definition of the allelic CAG repeat variant as given on page 5, lines 26-32 or lines 11 to 16, the claim is not rendered more clear. The definition of "alleles -3, -2, -1, 0, and 1" is an arbitrary one introduced by the applicant and therefore meaningless to the skilled person unless the full meaning is included in the claims. Finally, it would appear that claim 1 refers to human GT1 (hGT1). However, reference to Fig. 3 which shows a human and a mouse GT1 sequence, introduces some doubt whether the mouse sequence should be involved in the scope of claim 1 or not.

In view of the above, the dependent claims 2-12 are not clearly defined in the sense of Art. 6 PCT as well.

With respect to claim 8 it is further noted that it appears to be incomplete (A method to identify genes **part of or interacting with** a biochemical ...). Moreover, the claimed method is not defined by particular procedure steps to identify a gene which forms part of or which interacts with a biochemical pathway affected by the hGT1 gene. Screening of samples with probes or primers derived from the (wild type?) hGT1 sequence does not appear to result in the identification of the desired gene which interacts for example with the biochemical pathway but rather in the identification of allelic variants of hGT1 gene. As already mentioned above, it is not clear whether claim 9 refers to the (wildtype) hGT1 gene of claim 1 or to particular allelic CAG repeat variants thereof.

6). Section V

None of the available prior art discloses or suggests means and methods as described in the present application which therefore appears to conform with the

**INTERNATIONAL PRELIMINARY
EXAMINATION REPORT - SEPARATE SHEET**

International application No. PCT/CA98/00884

requirements of novelty and inventive step according to Art. 33(2), (3) PCT.

The embodiments of the invention in which an exclusive property or privilege is claimed are defined as follows:

1. A hGT1 gene containing transcribed polymorphic CAG repeat, which comprises a sequence as set forth in Fig. 3 and Figs. 4A-4C, wherein allelic variants of CAG repeat are selected from the group consisting of alleles -3, -2, -1, 0 and 1, and wherein said allelic variants are associated with schizophrenia, affective disorders, neurodevelopmental brain diseases or with phenotypic variability with respect to long term response to neuroleptic medication.

2. The gene of claim 1, wherein said affective disorder is manic depression.

3. A method for the prognosis of severity of schizophrenia of a patient, which comprises the steps of:

- a) obtaining a nucleic acid sample of said patient; and
- b) determining allelic variants of CAG repeat of the gene of claim 1, and wherein allelic variants shorter than allele 0 are indicative of non-severe schizophrenia.

4. A method for the identification of patient responding to neuroleptic medication, which comprises the steps of:

- a) obtaining a nucleic acid sample of said patient; and
- b) determining allelic variants of CAG repeat of the gene of claim 1, and wherein allelic

variants shorter than allele 0 are indicative of neuroleptic response.

5. The method of claim 4, wherein said shorter allelic variants have from about 171 to about 177 bp in length.

6. A non-human mammal model for the hGT1 gene of claim 1, whose germ cells and somatic cells are transformed and expresses at least one allelic variant of the hGT1 gene and wherein said allelic variant of the hGT1 being introduced into the mammal, or an ancestor of the mammal, at an embryonic stage.

7. A method for the screening of therapeutic agents for the prevention and/or treatment of schizophrenia, which comprises the steps of:

- a) administering said therapeutic agents to the non-human mammal of claim 6 or schizophrenia patients; and
- b) evaluating the prevention and/or treatment of development of schizophrenia in said mammal or said patients.

8. A method to identify genes part of or interacting with a biochemical pathway affected by hGT1 gene, which comprises the steps of:

- a) designing probes and/or primers using the hGT1 gene of claim 1 and screening psychiatric patients samples with said probes and/or primers; and
- b) evaluating the identified gene role in psychiatric patients.

9. A method of categorizing psychiatric patients according to their genotype to maximize response to treatment patients, which comprises the steps of:

- a) obtaining a nucleic acid sample of said patients; and
- b) determining allelic variants of CAG repeat of the gene of claim 1, wherein patients are categorized with respect to their allelic variants and wherein allelic variants shorter than allele 0 are indicative of neuroleptic response.

10. The use of the determination of allelic variants of CAG repeat of the gene of claim 1 for the identification of patient responding to neuroleptic medication, wherein allelic variants shorter than allele 0 are indicative of neuroleptic response.

11. The use of claim 10, wherein said shorter allelic variants have from about 171 to about 177 bp in length.

12. The use of the model of claim 6 for the screening of therapeutic agents for the manufacture of a medicament for prevention and/or treatment of schizophrenia.

AMENDED SHEET

- 25 -

The embodiments of the invention in which an exclusive property or privilege is claimed are defined as follows:

1. A hGT1 gene containing transcribed polymorphic CAG repeat, which comprises a sequence as set forth in Fig. 3 and Figs. 4A-4C.

2. The gene of claim 1, wherein allelic variants of CAG repeat are associated with schizophrenia, affective disorders, neurodevelopmental brain diseases or with phenotypic variability with respect to long term response to neuroleptic medication.

3. The gene of claim 2, wherein said affective disorder is manic depression.

4. A method for the prognosis of severity of schizophrenia of a patient, which comprises the steps of:

- a) obtaining a nucleic acid sample of said patient; and
- b) determining allelic variants of CAG repeat of the gene of claim 1, and wherein short allelic variants are indicative of non-severe schizophrenia.

5. A method for the identification of patient responding to neuroleptic medication, which comprises the steps of:

- a) obtaining a nucleic acid sample of said patient; and
- b) determining allelic variants of CAG repeat of the gene of claim 1, and wherein short allelic variants are indicative of neuroleptic response.

REPLACED by FIG 34

- 26 -

6. The method of claim 5, wherein said short allelic variants have from about 171 to about 177 bp in length.

7. A non-human mammal model for the hGT1 gene of claim 1, whose germ cells and somatic cells are modified to express at least one allelic variant of the hGT1 gene and wherein said allelic variant of the hGT1 being introduced into the mammal, or an ancestor of the mammal, at an embryonic stage.

8. A method for the screening of therapeutic agents for the prevention and/or treatment of schizophrenia, which comprises the steps of:

- a) administering said therapeutic agents to the non-human mammal of claim 7 or schizophrenia patients; and
- b) evaluating the prevention and/or treatment of development of schizophrenia in said mammal or said patients.

9. A method to identify genes part of or interacting with a biochemical pathway affected by hGT1 gene, which comprises the steps of:

- a) designing probes and/or primers using the hGT1 gene of claim 1 and screening psychiatric patients samples with said probes and/or primers; and
- b) evaluating the identified gene role in psychiatric patients.

10. A method of stratifying psychiatric patients based on the allelic variants of the hGT1 gene of claim 1 for clinical trials purposes, which comprises:

- 27 -

- a) obtaining a nucleic acid sample of said patients; and
- b) determining allelic variants of CAG repeat of the gene of claim 1, wherein patients are stratified with respect to their allelic variants and wherein short allelic variants are indicative of neuroleptic response.

PCT

INTERNATIONAL SEARCH REPORT

(PCT Article 18 and Rules 43 and 44)

Applicant's or agent's file reference 1770-182PCT	FOR FURTHER ACTION see Notification of Transmittal of International Search Report (Form PCT/ISA/220) as well as, where applicable, item 5 below.	
International application No. PCT/CA 98/ 00884	International filing date (day/month/year) 18/09/1998	(Earliest) Priority Date (day/month/year) 19/09/1997
Applicant McGILL UNIVERSITY et al.		

This International Search Report has been prepared by this International Searching Authority and is transmitted to the applicant according to Article 18. A copy is being transmitted to the International Bureau.

This International Search Report consists of a total of 3 sheets.

☒ It is also accompanied by a copy of each prior art document cited in this report.

1. ☐ Certain claims were found unsearchable (see Box I).

2. ☐ Unity of invention is lacking (see Box II).

3. ☒ The international application contains disclosure of a **nucleotide and/or amino acid sequence listing** and the international search was carried out on the basis of the sequence listing

☐ filed with the international application.

☒ furnished by the applicant separately from the international application.

☐ but not accompanied by a statement to the effect that it did not include matter going beyond the disclosure in the international application as filed.

☐ Transcribed by this Authority

4. With regard to the **title**, ☐ the text is approved as submitted by the applicant

☒ the text has been established by this Authority to read as follows:

POLYMORPHIC CAG REPEAT-CONTAINING GENE AND USES THEREOF

5. With regard to the **abstract**,

☒ the text is approved as submitted by the applicant

☐ the text has been established, according to Rule 38.2(b), by this Authority as it appears in Box III. The applicant may, within one month from the date of mailing of this International Search Report, submit comments to this Authority.

6. The figure of the **drawings** to be published with the abstract is:

Figure No. 3 ☐ as suggested by the applicant.

☐ None of the figures.

☒ because the applicant failed to suggest a figure.

☐ because this figure better characterizes the invention.

INTERNATIONAL SEARCH REPORT

National Application No.

T/CA 98/00884

A. CLASSIFICATION OF SUBJECT MATTER

IPC 6 C12N15/00 C12Q1/68 C07K14/47 A01K67/027

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 6 C07K C12Q

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	IMAI, Y. ET AL.: "Cloning of a retinoic acid-induced gene, GT1, in the embryonal carcinoma cell line P19: neuron-specific expression in the mouse brain" MOLECULAR BRAIN RESEARCH, XP002093328 see the whole document ---	1
A	ROBITAILLE, Y. ET AL.: "The neuropathology of CAG repeat diseases: review and update of genetic and molecular features" BRAIN PATHOLOGY, vol. 7, no. 3, July 1997, pages 901-926, XP002093329 see the whole document --- -/--	2



Further documents are listed in the continuation of box C.



Patent family members are listed in annex.

* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

"&" document member of the same patent family

Date of the actual completion of the international search

12 February 1999

Date of mailing of the international search report

25/02/1999

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Chambonnet, F

INTERNATIONAL SEARCH REPORT

International Application No

PCT/CA 98/00884

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category °	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	WO 95 01437 A (UNIV MINNESOTA) 12 January 1995 see the whole document ---	2-10
A	WO 97 18825 A (UNIV BRITISH COLUMBIA ;KALCHMAN MICHAEL (CA); HAYDEN MICHAEL R (CA) 29 May 1997 see the whole document ---	2-10
A	MACIEL, P. ET AL.: "Correlation between CAG repeat length and clinical features in Machado-Joseph Disease" AMERICAN JOURNAL OF HUMAN GENETICS, vol. 57, no. 1, July 1995, pages 54-61, XP002093330 ---	1,2,4,10
A	JOOPER, R. ET AL.: "Apolipoprotein E genotype in Schizophrenia" AMERICAN JOURNAL OF MEDICAL GENETICS (NEUROPSYCHIATRIC GENETICS), vol. 67, no. 2, 9 April 1996, page 235 XP002093331 see the whole document ---	2
P,X	PHILIBERT, R.A. ET AL.: "The characterization and sequence analysis of thirty CTG-repeat containing genomic cosmid clones" EUROPEAN JOURNAL OF HUMAN GENETICS, vol. 6, no. 1, January 1998, pages 89-94, XP002093332 see the whole document ---	1
T	TURECKI, G. ET AL.: "Schizophrenia and chromosome 6p" AMERICAN JOURNAL OF MEDICAL GENETICS, vol. 74, no. 2, 1997, pages 195-198, XP002093333 see the whole document -----	2

INTERNATIONAL SEARCH REPORT

Information on patent family members

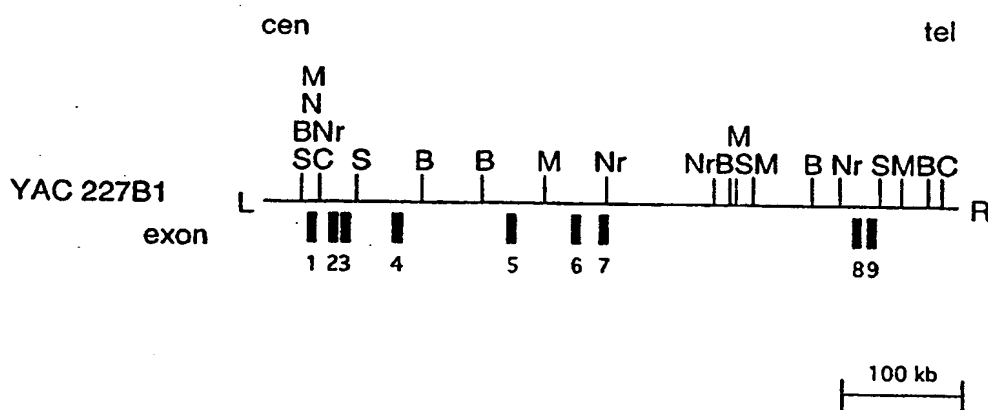
International Application No

PCT/CA 98/00884

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
WO 9501437 A	12-01-1995	US 5834183 A	10-11-1998
		CA 2166117 A	12-01-1995
		EP 0707647 A	24-04-1996
		JP 9501049 T	04-02-1998
		US 5741645 A	21-04-1998
WO 9718825 A	29-05-1997	CA 2238075 A	29-05-1997
		EP 0873132 A	28-10-1998

(51) International Patent Classification ⁶ : C12N 15/12, C12Q 1/68, C07K 14/47, 16/18, G01N 33/577		A2	(11) International Publication Number: WO 95/01437
			(43) International Publication Date: 12 January 1995 (12.01.95)
(21) International Application Number: PCT/US94/07336		(81) Designated States: CA, JP, European patent (AT, BE, CH, DE, DK, ES, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE).	
(22) International Filing Date: 29 June 1994 (29.06.94)		Published Without international search report and to be republished upon receipt of that report.	
(30) Priority Data: 08/084,365 29 June 1993 (29.06.93) US 08/267,803 28 June 1994 (28.06.94) US			
(71) Applicant: REGENTS OF THE UNIVERSITY OF MINNESOTA [US/US]; Morrill Hall, 100 Church Street, S.E., Minneapolis, MN 55455 (US).			
(72) Inventors: ORR, Harry, T.; 5133 Luverne Avenue South, Minneapolis, MN 55419 (US). CHUNG, Ming-yi; 425 13th Avenue South, Apartment 1405, Minneapolis, MN 55414 (US). ZOGHBI, Huda, Y.; 5801 Charlotte, Houston, TX 77005 (US).			
(74) Agent: RAASCH, Kevin, W.; Schwegman, Lundberg & Woessner, 3500 IDS Center, 80 South Eighth Street, Minneapolis, MN 55402 (US).			

(54) Title: GENE SEQUENCE FOR SPINOCEREBELLAR ATAXIA TYPE 1 AND METHOD FOR DIAGNOSIS



The present invention provides an isolated DNA molecule of the autosomal dominant spinocerebellar ataxia type 1 gene, which is located within the short arm of chromosome 6. This isolated DNA molecule is preferably located within a 3.36 kb *EcoRI* fragment, i.e., an *EcoRI* fragment containing about 3360 base pairs, of the SCA1 gene. The isolated sequences contain a CAG repeat region. The number of CAG trinucleotide repeats (n) is ≤ 36 , preferably $n = 19-36$, for normal individuals. For an affected individual $n > 36$, preferably $n \geq 43$.

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AT	Austria	GB	United Kingdom	MR	Mauritania
AU	Australia	GE	Georgia	MW	Malawi
BB	Barbados	GN	Guinea	NE	Niger
BE	Belgium	GR	Greece	NL	Netherlands
BF	Burkina Faso	HU	Hungary	NO	Norway
BG	Bulgaria	IE	Ireland	NZ	New Zealand
BJ	Benin	IT	Italy	PL	Poland
BR	Brazil	JP	Japan	PT	Portugal
BY	Belarus	KE	Kenya	RO	Romania
CA	Canada	KG	Kyrgyzstan	RU	Russian Federation
CF	Central African Republic	KP	Democratic People's Republic of Korea	SD	Sudan
CG	Congo	KR	Republic of Korea	SE	Sweden
CH	Switzerland	KZ	Kazakhstan	SI	Slovenia
CI	Côte d'Ivoire	LI	Liechtenstein	SK	Slovakia
CM	Cameroon	LK	Sri Lanka	SN	Senegal
CN	China	LU	Luxembourg	TD	Chad
CS	Czechoslovakia	LV	Latvia	TG	Togo
CZ	Czech Republic	MC	Monaco	TJ	Tajikistan
DE	Germany	MD	Republic of Moldova	TT	Trinidad and Tobago
DK	Denmark	MG	Madagascar	UA	Ukraine
ES	Spain	ML	Mali	US	United States of America
FI	Finland	MN	Mongolia	UZ	Uzbekistan
FR	France			VN	Viet Nam
GA	Gabon				

- 1 -

GENE SEQUENCE FOR SPINOCEREBELLAR ATAXIA TYPE 1
AND METHOD FOR DIAGNOSIS

5

Statement of Government Rights

The present invention was made with government support under Grant Nos. NS 22920 and 27699, awarded by the National Institutes of Health. The
10 Government has certain rights in this invention.

Background of the Invention

The spinocerebellar ataxias are a heterogeneous group of degenerative neurological disorders with variable clinical features resulting from
15 degeneration of the cerebellum, brain stem, and spinocerebellar tracts. The clinical symptoms include ataxia, dysarthria, ophthalmoparesis, and variable degrees of motor weakness. The symptoms usually begin during the third or fourth decade of life, however, juvenile onset has been identified. Typically, the disease worsens gradually, often resulting in complete disability and death 10-20 years after the
20 onset of symptoms. Individuals with juvenile onset spinocerebellar ataxias, however, typically have more rapid progression of the phenotype than the late onset cases. A method for diagnosing spinocerebellar ataxias would provide a significant step toward its treatment.

Spinocerebellar ataxia type 1 (SCA1) is an autosomal dominant
25 disorder which is genetically linked to the short arm of chromosome 6 based on linkage to the human major histocompatibility complex (HLA). See, for example, H. Yakura et al., N. Engl. J. Med., 291, 154-155 (1974); and J.F. Jackson et al., N. Engl. J. Med., 296, 1138-1141 (1977). SCA1 has been shown to be tightly linked to the marker D6S89 on the short arm of chromosome 6, telomeric to HLA. See, for
30 example, L.P.W. Ranum et al., Am. J. Hum. Genet., 49, 31-41 (1991); and H.Y. Zoghbi et al., Am. J. Hum. Genet., 49, 23-30 (1991). Recently, two families with dominantly inherited ataxia failed to show detectable linkage with HLA markers but were found to have SCA1 when studied for linkage to D6S89, demonstrating the superiority of the latter marker for study of ataxia families. See, for example, B.J.B.
35 Keats et al., Am. J. Hum. Genet., 49, 972-977 (1991). The identification and cloning of the SCA1 gene could provide methods of detection that would be extremely valuable for both family counseling and planning medical treatment.

Summary of the Invention

The present invention is directed to a portion of an isolated 1.2-Mb region of DNA from the short arm of chromosome 6 containing a highly polymorphic CAG repeat region in the SCA1 gene. This CAG repeat region is unstable (i.e., highly variable within a population) and is expanded in individuals with the autosomal dominant neurodegenerative disorder spinocerebellar ataxia type 1 (i.e., affected individuals generally have more than 36 CAG repeats). Southern and PCR analyses of the CAG repeat region demonstrate correlation between the size of the expanded repeat region and the age-of-onset of the disorder (with larger alleles, i.e., more repeat units, occurring in juvenile cases), and severity of the disorder (with larger alleles, i.e., more repeat units, occurring in the more severe cases).

Specifically, the present invention provides a nucleic acid molecule containing a CAG repeat region of an isolated autosomal dominant spinocerebellar ataxia type 1 gene (herein referred to as "SCA1"), which is located within the short arm of chromosome 6. The SCA1 gene contains a region that encodes a protein herein referred to as "ataxin-1." The nucleic acid molecule of the present invention can be a single or a double-stranded polynucleotide. It can be genomic DNA, cDNA, or mRNA of any size as long as it includes the CAG repeat region of an isolated SCA1 gene. Preferably, the nucleic acid molecule includes the SCA1 coding region and is of about 2.4-11 kb in length. It can be the entire SCA1 gene (whether genomic DNA or a transcript thereof) or any fragment thereof that contains the CAG region of the gene. One such fragment is an *EcoRI* fragment of the SCA1 gene, i.e., a fragment obtained through digestion with *EcoRI* endonuclease restriction enzyme, containing about 3360 base pairs having therein a polymorphic CAG repeat region. By polymorphic CAG repeat region it is meant that there are repeating CAG trinucleotides in this portion of the gene that can vary in the number of CAG trinucleotides. The number of trinucleotide repeats can vary from as few as 19, for example, to as many as 81, for example, and larger.

For a normal individual, $n \leq 36$ in the $(CAG)_n$ region, i.e., $n = 2-36$, and typically $n = 19-36$. This region in a normal allele of the SCA1 gene is optionally interrupted with CAT trinucleotides. Typically, there are no more than about 3 CAT trinucleotides, either individually or in combination, within any

-3-

(CAG)_n region. The (CAG)_n region of this isolated sequence is unstable, i.e., highly variable within a population, and larger, i.e., expanded, in individuals who have symptoms of the disease, or who are likely to develop symptoms of the disease. For an affected individual, i.e., an individual with an affected allele of the SCA1 gene, n > 36 in the (CAG)_n region, and typically n ≥ 43. One isolated DNA molecule of the SCA1 gene is about 3360 base pairs in length as shown in Figure 1. The sequences of a portion of the *Eco*RI fragment within the SCA1 gene of several affected individuals is shown in Figure 2. The entire 10,660 nucleotides of the SCA1 gene transcript are shown in Figure 15 (the entire SCA1 gene spans about 450 kb of genomic DNA).

The present invention is also directed to isolated oligonucleotides, particularly primers for use in PCR techniques and probes for diagnosing the neurodegenerative disorder SCA1. The oligonucleotides have at least about 11 nucleotides and hybridize to a nucleic acid molecule containing a CAG repeat region of an isolated SCA1 gene. The hybridization can occur to any portion of a nucleic acid molecule containing a CAG repeat region of the SCA1 gene. Preferably, the oligonucleotides hybridize to a 3.36 kb *Eco*RI fragment of an SCA1 gene having a CAG repeat region. Alternatively stated, each oligonucleotide is substantially complementary (having greater than 65% homology) to a nucleotide sequence having a CAG repeat region, i.e., a (CAG)_n region, preferably to a 3.36-kb *Eco*RI fragment of the SCA1 gene. If the oligonucleotide is a primer the molecule preferably contains at least about 16 nucleotides and no more than about 35 nucleotides. Furthermore, preferred primers are chosen such that they produce a primed product of about 70-350 base pairs, preferably about 100-300 base pairs. More preferably, the primers are chosen such that nucleotide sequence is complementary to a portion of a strand of an affected or a normal allele within about 150 nucleotides on either side of the (CAG)_n region, including directly adjacent to the (CAG)_n region. Most preferably, the primer is selected from the group consisting of CCGGAGCCCTGCTGAGGT (CAG-a), CCAGACGCCGGGACAC (CAG-b), AACTGGAAATGTGGACGTAC (Rep-1), CAACATGGGCAGTCTGAG (Rep-2), CCACCACTCCATCCCAGC (GCT-435), TGCTGGGCTGGTGGGGGG (GCT-214), CTCTCGGCTTTCTTGGTG (Pre-1), and GTACGTCCACATTTCCAGTT (Pre-2). These primers substantially correspond to those shown in Figure 3.

-4-

They can be used in any combination for sequencing or producing amplified nucleic acid molecules, e.g., DNA molecules, using various PCR techniques. Preferably, for amplification of the DNA molecule characteristic of the SCA1 disorder, Rep-1 and Rep-2 is the primer pair used. As used herein, the term

5 "amplified DNA molecule" refers to DNA molecules that are copies of a portion of DNA and its complementary sequence. The copies correspond in nucleotide sequence to the original DNA sequence and its complementary sequence. The term "complement", as used herein, refers to a DNA sequence that is complementary (having greater than 65% homology) to a specified DNA sequence. The term

10 "primer pair", as used herein, means a set of primers including a 5' upstream primer that hybridizes with the 5' end of the DNA molecule to be amplified and a 3' downstream primer that hybridizes with the complement of the 3' end of the molecule to be amplified.

Using the primers of the present invention, PCR technology can be

15 used in the diagnosis of the neurological disorder SCA1 by detecting a region of greater than about 36 CAG repeating trinucleotides, preferably at least 43 repeating CAG trinucleotides. Generally, this involves treating separate complementary strands of the DNA molecule containing a region of repeating CAG codons with a molar excess of two oligonucleotide primers, extending the primers to form

20 complementary primer extension products which act as templates for synthesizing the desired molecule containing the CAG repeating units, and detecting the molecule so amplified.

An oligonucleotide that can be used as a gene probe for identifying a nucleic acid molecule, e.g., a DNA molecule, containing a CAG repeat region of the

25 SCA1 gene is also provided. The gene probe can be used for distinguishing between the normal and the larger affected alleles of the SCA1 gene. The gene probe can be a portion of a nucleotide sequence of the SCA1 gene itself (e.g., a 3.36-kb *EcoRI* fragment or portion thereof), complementary to it, or hybridizable to it or the complement. It is of a size suitable for forming a stable duplex, i.e., having

30 at least about 11 nucleotides, preferably having at least about 15 nucleotides, more preferably having at least about 100 nucleotides (for effective Southern blotting), and most preferably having at least about 200 nucleotides. The probe can contain any portion of the (CAG)_n region, although this is not a requirement. It is desirable, however, for the probe to contain a portion of the nucleic acid molecule on either

-5-

side of the (CAG)_n region. There is generally no maximum size limitation for such probes. In fact, the entire SCA1 gene could be a probe.

The gene probe of the present invention is useable in a method of diagnosing a patient for SCA1. A particularly preferred method of diagnosis involves detecting the presence of a DNA molecule containing a CAG repeat region of the SCA1 gene. Specifically, the method includes the steps of digesting genomic DNA with a restriction endonuclease to obtain DNA fragments; preferably, separating the fragments by size using gel electrophoresis; probing said DNA fragments under hybridizing conditions with a detectably labeled gene probe that hybridizes to a nucleic acid molecule containing a CAG repeat region of an isolated SCA1 gene; detecting probe DNA which has hybridized to said DNA fragments; and analyzing the DNA fragments for a (CAG)_n region characteristic of the normal or affected forms of the SCA1 gene.

The present invention also provides a protein (or portions thereof) encoded by the SCA1 gene and antibodies (polyclonal or monoclonal) produced from the protein or portions thereof. The antibodies can be used in methods of isolating antigenic protein expressed by the SCA1 gene. For example, they can be added to a biological sample containing the antigenic protein to form an antibody-antigen complex, which can be isolated from the sample and exposed to amino acid sequencing of the antigenic protein. This can be done while the protein is still complexed with the antibody.

Thus, the present invention provides methods to determine the presence or absence of an affected form of the SCA1 gene, which can be based on RNA- or DNA-based detection methods (preferably, the methods involve isolating and analyzing genomic DNA) or on protein-based detection methods. These methods include, for example, PCR-based methods, direct nucleic acid sequencing, measuring expression of the SCA1 gene by measuring the amount of mRNA expressed or by measuring the amount of ataxin-1 protein expressed. The methods of the present invention also include determining the size of the repeat region of the nucleic acid or amino acid molecules.

As used herein, the term "isolated (and purified)" means that the nucleic acid molecule, gene, or oligonucleotide is essentially free from the remainder of the human genome and associated cellular or other impurities. This does not mean that the product has to have been extracted from the human genome;

-6-

rather, the product could be a synthetic or cloned product for example. As used herein, the term "nucleic acid molecule" means any single or double-stranded RNA or DNA molecule, such as mRNA, cDNA, and genomic DNA.

As used herein, the term "SCA1 gene" means the
5 deoxyribopolynucleotide located within the short arm of chromosome 6 between markers D6S89 and D6S274 of about 450 kb (10.5-11 kb transcript) containing an unstable CAG repeat region. This term, therefore, refers to numerous unique genes that are substantially the same except for the content of the CAG repeat region. A representative example of the SCA1 gene transcript for a normal individual is shown
10 in Figure 15. Included within the scope of this term is any ribo- or deoxyribo-polynucleotide containing zero, one or more nucleotide substitutions that also encodes the protein ataxin-1. Included in the term "SCA1 gene" is any polynucleotide as described in the previous sentence that has different numbers of CAG and/or CAT repeats in the polymorphic CAG repeat region. It is understood
15 also that the term "SCA1 gene" includes both the polypeptide-encoding region and the regions that encode the 5' and 3' untranslated segments of the mRNA for SCA1. Although the SCA1 gene described herein is described in terms of the human genome, it is envisioned that other mammals, e.g., mice, may also have a very similar gene containing a CAG repeat region that could be used to produce
20 oligonucleotides, for example, that are useful in diagnosing the SCA1 disorder in humans.

As used herein, the term "ataxin-1" means the gene product of the SCA1 gene, i.e., protein encoded by the open reading frame of the SCA1 gene and any protein substantially equivalent thereto, including all proteins of different
25 lengths (e.g., 20-90 kD, preferably 60-90 kD) encoded by said open reading frame which start at each in-frame ATG translation start site. The term "ataxin-1" further includes all proteins with essentially the same N-terminal and C-terminal sequences but different numbers of glutamine (Q) and/or histidine (H) repeats (primarily glutamine repeats) in the polymorphic repeat region.

30 As used herein, the term "polymorphic CAG repeat region" or simply "CAG repeat region" means that region of the SCA1 gene that encodes a string of polyglutamate residues that varies in number from individual allele to individual allele, and which can range in number from 2 to 80 or more. Moreover, the polymorphic CAG repeat regions can contain CAT (encoding histidine) in place of

-7-

CAG, although CAT is much less common than CAG in this region. It is to be understood that when referring to nucleic acid molecules containing the CAG repeat region, this includes RNA molecules containing the corresponding GUC repeat region.

5 As used herein, an "affected" gene refers to the allele of the SCA1 gene that, when present in an individual, is the cause of spinocerebellar ataxia type 1, and an "affected" individual has the symptoms of autosomal dominant spinocerebellar ataxia type 1. Individuals with only "normal" SCA1 genes, do not possess the symptoms of SCA1. The term "allele" means a genetic variation
10 associated with a coding region; that is, an alternative form of the gene.

As used herein, "hybridizes" means that the oligonucleotide forms a noncovalent interaction with the stringency target nucleic acid molecule under standard conditions. The hybridizing oligonucleotide may contain nonhybridizing nucleotides that do not interfere with forming the noncovalent interaction, e.g., a
15 restriction enzyme recognition site to facilitate cloning.

Brief Description of the Drawings

Figure 1. Sequence of the 3.36 kb *Eco*RI fragment of the normal SCA1 gene located within the short arm of chromosome 6. It is within this
20 fragment that mutations occur in the CAG repeat region which are associated with autosomal dominant spinocerebellar ataxia type 1.

Figure 2. Sequence information for five affected individuals in the CAG repeat region, i.e., the CAG trinucleotide repeat, and its flanking regions of the SCA1 gene located within a short arm of chromosome 6.

25 **Figure 3.** Sequence of the CAG trinucleotide repeat and its flanking regions. About 500 nucleotides in a single strand of DNA of the 3.36 kb *Eco*RI fragment of the SCA1 gene shown in Figure 1 is represented. The locations of PCR primers are shown by solid lines with arrowheads.

Figure 4. Summary of SCA1 recombination events that led to the
30 precise mapping of the SCA1 locus. Recombinant disease-carrying chromosomes are shown for the markers shown above. A schematic diagram of the relevant region of 6p22 (not drawn to scale) is shown at the top of the figure. Families are coded as follows: TX = Houston, MN = Minnesota, MI = Michigan, IT = Italy. Each recombination event is given a number following the family code.

-8-

Figure 5. Regional localization of 6p22-p23 STSs by PCR analysis of radiation reduced hybrids. Three panels (a-c) demonstrate the regional localization of D6S274, D6S288, and AM10GA. In each panel PCR amplification results are shown for genomic DNA, the I-7 cell line which retains 6p, the radiation reduced hybrids R17, R72, R86, and R54, and RJK88 hamster DNA. A blank control (c) is shown for every panel. R86 has been previously shown to retain D6S89; R17 and R72 are known to contain D6S88 and D6S108, two DNA markers which map centromeric to D6S89. An amplification product is seen in I-7, R17, R72, and R86 for D6S274 and D6S288, whereas the amplification product for AM10GA is only seen in I-7 and R86 confirming that D6S274 and D6S288 map centromeric to AM10GA and D6S89.

Figure 6. A schematic diagram of 6p22-p23 region showing the new markers and the YAC contig. At the bottom of the diagram, the radiation hybrid reduced panel used for regional mapping is shown. YAC clones are represented as dark lines, open segments indicate a noncontiguous region of DNA. The discontinuity shown in YAC clone 351B10 indicate that this YAC has an internal deletion. All of the ends of the YAC clones that were isolated are designated by an "L" for the left end or an "R" for the right end.

Figure 7. Genotypic data for 6p22-p23 dinucleotide repeat markers are shown for a reduced pedigree from the MN-SCA1 kindred. This figure summarizes a second recombination event that led to the precise mapping of the SCA1 locus.

Figure 8. Long-range restriction maps of YACs, 227B1, 60H7, 195B5, A250D5, and 379C2. YACs 351B10, 172B5, 172B5, and 168F1 were also used in the restriction analysis (data not shown). The restriction sites are marked as N, *NotI*; B, *BssHII*; Nr, *NruI*; M, *MluI*; S, *SacII*, and Sa, *SalI*. A summary map of the SCA1 gene region with the position of the DNA markers used as probes (boxes) is shown. The centromere-telomere orientation is indicated by cen/tel respectively.

Figure 9. Physical map of the SCA1 region. The positions of various genetic markers and sequence tagged sites (STSs) relative to the overlapping YAC clones are shown. AM10 and FLB1 are STSs developed using a radiation reduced hybrid retaining chromosome 6p22-p23, A205D5-L and 195B5-L are STSs from insert termini of YACs A250D5 and 195B5. D6S89, D6S109, D6S288 and D6S274, and AM10-GA are dinucleotide repeat markers used in the genetic analysis

of SCA1 families. The SCA1 candidate region is flanked by the D6S274 and D6S89 markers which identify the closest recombination events. The YAC clones shown here are indicated by the cross-hatched markings. YAC 172B5 has two non-contiguous segments of DNA as indicated by the open bar for the non-6p segment. The YACs are designated according to St. Louis and CEPH libraries. The position of the cosmid contig (C) which contains the overlapping cosmids which are (CAG)_n positive is indicated by a solid black bar. The overlap between the YACs was determined by long-range restriction analysis. Orientation is indicated as centromeric (Cen) and telomeric (Tel).

Figure 10. Southern blot analysis of leukocyte DNA using the 3.36-kb *EcoRI* fragment which contains the repeat as a probe. **Figure 10a:** *TaqI*-digested DNA from a TX-SCA1 kindred. The unaffected spouse has a single fragment at 2830-bp. The affected individual with onset at 25 years of age has the 2830-bp fragment as well as a 2930-bp fragment. The affected child with onset at 4 years inherited the normal 2830-bp from her mother, and has a new fragment of 3000-bp not seen in either parent. **Figure 10b:** *TaqI*-digested DNA from individuals from a MN-SCA1 kindred. The unaffected spouse and the unaffected sibling have a 2830-bp fragment. The two affected brothers have the 2830-bp fragment as well as an expanded fragment of 2900-bp in the sib with onset at 25 years and 2970-bp in the sib with onset at 9 years. **Figure 10c:** *BstNI*-digested DNA from the TX-SCA1 kindred. Lanes 1-3 are from the same kindred depicted in (A). The normal fragment size is 530-bp, in individuals with onset at 25-30 years (lanes 1 and 4) the fragment expands to 610-bp. In the individual with onset at 15 years of age (lane 7) the fragment size is 640-bp, and in the individual with onset at 4 years (lane 3) the fragment size is 680-bp. The DNA in lane 5 is from a 14 year old child who is asymptomatic.

Figure 11. Analysis of the PCR-amplified products containing the trinucleotide repeat tract in normal and SCA1 individuals. The CAG-a/CAG-b primer pair was used in panel (a) whereas the Rep-1/Rep-2 primer pair was used in panel (b). The individuals in lanes 1, 2 and 3 in panel (a) are brothers. The range for the normal (NL) and expanded (EXP) (CAG)_n repeat units is indicated.

Figure 12. A scatter plot for the age-at-onset in years versus the number of the (CAG)_n repeat units is shown to demonstrate the correlation between the age-at-onset and the size of the expansion. A linear correlation coefficient of

-10-

-0.845 was obtained. In addition a curvilinear correlation coefficient was calculated given the non-linear pattern of the plot. The curvilinear correlation coefficient is -0.936.

Figure 13. Schematic representation of the SCA1 cDNA contig. A subset of overlapping phage cDNA clones (black bars) and 5'-RACE-PCR product (R1) spanning 10.66 kb of the SCA1 transcript is shown. cDNA clone 31-5 contains the entire coding region for the SCA1 gene product, ataxin-1. On top, a schematic shows the structure of the SCA1 transcript; the sizes of the coding region (rectangle) as well as the 5'UTR and the 3'UTR (thin lines) are indicated. The position of the CAG repeat within the coding region is also shown. An asterisk indicates the clones used as probes to screen the cDNA libraries. At the bottom the positions of *Bam*HI (B), *Hind*III (H), and *Taq*I (T) restriction sites are shown.

Figure 14. Northern blot analysis of the SCA1 gene using RNAs from multiple human tissues. The panel on the left is probed with a PCR product from a portion of the coding region (bp 2460 to bp 3432). The panel on the right is hybridized with the 3J cDNA clone from the 3'UTR. An ~11 kb transcript is detected in RNAs from all tissues using both probes as well as the cDNA clones 31-5 and 8-8, both of which contain the CAG repeat (Figure 13).

Figure 15. The sequence of the SCA1 transcript. The sequences of primers 9b, 5F and 5R (bp 129-147, bp 173-191 and bp 538-518 respectively in the 5' to 3' orientation) are underlined. The protein sequence encoded by the DNA is shown below the DNA sequence. The CAG repeat region is from about bp 1524 to about bp 1613.

Figure 16. a. The structure of the SCA1 transcript and the various splice variants. The schematic on top represents the nine exons (not drawn to scale) and their respective sizes. The stippled areas indicate the coding region. The structure of five cDNA clones representing different splice variants of the SCA1 transcript are also shown. Clones 8-8 and 8-9b are phage clones, RT-PCR1 and RT-PCR2 are two clones obtained by RT-PCR carried out on cerebellar poly-(A)⁺ RNA using the primers 9b and 5R (Figure 15). Only 30 bp of exon 1 were present in clone 8-9b and RT-PCR products as indicated by the broken line in the rectangles. **b.** Detection of alternative splicing of the SCA1 transcript in cerebellar poly-(A)⁺ RNA (CBL RNA). RT-PCR analysis was carried out using two sets of primers: 9b-5R and 5F-5R. PCR products of the expected size were detected in

-11-

CBL RNA in the presence of reverse transcriptase (+RT) with both pairs of primers. Using the 9b-5R pair at least two larger PCR products were also detected. Using the 5F-5R pair for RT-PCR at annealing $T < 60^\circ$, some faint bands in the same size range as those seen using the 9b-5R primer pair were also seen. 8-8 and 8-9b are the
 5 phage clones used as positive controls. The sizes of the relevant bands of the molecular weight marker (FX174 cut with *HaeIII*) are indicated on the left.

Figure 17. Intron-exon boundaries of the SCA1 gene. Splice acceptor and splice donor sites are indicated in bold letters. The numbers at the beginning and the end of each exon refer to the position in the composite sequence
 10 of SCA1 in Figure 15. Uppercase letters indicate exon sequences, lowercase letters indicate intron sequences. Y= pyrimidine; R= purine; N= undefined.

Figure 18. Genomic structure of the SCA1 gene. The nine exons of the SCA1 gene (solid rectangles not drawn to scale) were localized based on the restriction map of the SCA1 region by Southern analysis using rare cutter DNA
 15 digests from several YAC clones. A representative map using YAC clone 227B1, which encompasses the SCA1 gene, is shown. The restriction map of this YAC has been confirmed by analysis of four overlapping YAC clones in the region. The centromere-telomere orientation is indicated by CEN-TEL, respectively. L= left YAC end; R= right YAC end; B= *BssHIII*; C= *CspI*; M= *MluI*; N= *NotI*; Nr= *NruI*;
 20 S= *SacII*.

Figure 19. Analysis of expression of the expanded SCA1 allele. RT-PCR was carried out on lymphoblast poly-(A)⁺RNA from one unaffected individual (lane 1) and four SCA1 patients (lanes 2 through 5) using primers Rep1 and Rep2. This analysis shows that both the normal and the expanded SCA1 alleles
 25 are transcribed. The number of the repeat units for each allele is indicated below each lane; lane 6 is the RT minus control.

Figure 20. Distributions of CAG repeat lengths from unaffected control individuals and from SCA1 alleles. Normal alleles range in size from 19 to 36 repeat units while disease alleles contain from 42 to 81 repeats.

30

-12-

Detailed Description

Substantial efforts have been made to localize the SCA1 gene using genetic and physical mapping methods. Genetically, SCA1 is flanked on the centromeric side by D6S88 at a *rE*combination fraction of approximately 0.08 (based on marker-marker distances using the Centre d'Etude du Polymorphisme Humain (CEPH) reference families) and on the telomeric side by F13A at a recombination fraction of 0.19. See, L.P.W. Ranum et al., Am. J. Hum. Genet., 49, 31-41 (1991). Both markers are quite distant and are not practical for use in efforts aimed at cloning the SCA1 gene. The D6S89 marker maps closer to the SCA1 gene.

To localize SCA1 more precisely, five dinucleotide polymorphisms near D6S89 have been identified. A new marker, AM10GA, demonstrates no recombination with SCA1. Linkage analysis and analysis of recombination events confirm that SCA1 maps centromeric to D6S89 with D6S109 as the other flanking marker at the centromeric end and establishes the following order: centromere-D6S109-AM10GA/SCA1-D6S89-LR40-D6S202-telomere. The genetic distance between the two flanking markers D6S109 and D6S89 is about 6.7 cM based on linkage analysis using 40 reference families from the Centre d'Etude du Polymorphisme Humain (CEPH).

A. SCA1 Gene and Method of Diagnosis

The size of the candidate region on the short arm of chromosome 6 containing the SCA1 locus is about 1.2 Mb, and is flanked by D6S274 to the centromeric side and D6S89 to the telomeric side. The SCA1 gene spans 450 kb of genomic DNA and is organized in nine exons (Figure 15 is representative of the SCA1 gene from a normal individual). The SCA1 transcript (i.e., mRNA or cDNA clone) is about 10.6-11 kb. The gene is transcribed in both normal and affected SCA1 alleles. The structure of the gene is unusual in that it contains seven exons in the 5'-untranslated region, two large exons (2080 bp and 7805 bp) which contain a 2448-bp coding region, and a 7277 bp 3'-untranslated region. The first four non-coding exons undergo extensive alternative splicing in several tissues.

The gene for SCA1 contains a highly polymorphic CAG repeat that is located within a 3.36-kb fragment produced by digestion of the candidate region with the restriction enzyme, *EcoRI*. The CAG repeat region preferably lies within

-13-

the coding region and codes for polyglutamine. This region of CAG repeating sequences is unstable and expanded in individuals with SCA1. Southern and PCR analyses of the (CAG)_n repeat demonstrate a correlation between the size of the repeat expansion and the age-at-onset of SCA1 and severity of the disorder. That is, individuals with more repeat units (or longer repeat tracts) tend to have both an early age of onset and a more severe disease course. These results demonstrate that SCA1, like fragile X syndrome, myotonic dystrophy, X-linked spinobulbar muscular atrophy, and Huntington disease, displays a mutational mechanism involving expansion of an unstable trinucleotide repeat.

The identification of a trinucleotide repeat expansion associated with SCA1 allows for improved diagnosis of the disease. Thus, in addition to being directed to the gene for SCA1 and the protein encoded thereby, the present invention also relates to methods of diagnosing SCA1. These diagnostic methods can involve any known method for detecting a specific fragment of DNA. These methods can include direct detection of the DNA or indirect through detection of RNA or proteins, for example. For example, Southern or Northern blotting hybridization techniques using labeled probes can be used. Alternatively, PCR techniques can be used with novel primers that amplify the CAG repeating region of the *EcoRI* fragment. Nucleic acid sequencing can also be used as a direct method of determining the number of CAG repeats.

For example, DNA probes can be used for identifying DNA segments of the affected allele of the SCA1 gene. DNA probes are segments of labeled, single-stranded DNA which will hybridize, or noncovalently bind, with complementary single-stranded DNA derived from the gene sought to be identified. The probe can be labeled with any suitable label known to those skilled in the art, including radioactive and nonradioactive labels. Typical radioactive labels include ³²P, ¹²⁵I, ³⁵S, and the like. Nonradioactive labels include, for example, ligands such as biotin or digoxigenin as well as enzymes such as phosphatase or peroxidases, or the various chemiluminescers such as luciferin, or fluorescent compounds like fluorescein and its derivatives. The probe may also be labeled at both ends with different types of labels for ease of separation, as, for example, by using an isotopic label at one end and a biotin label at the other end.

Using DNA probe analysis, the target DNA can be derived by the enzymatic digestion, fractionation, and denaturation of genomic DNA to yield a

-14-

complex mixture incorporating the DNA from many different genes, including DNA from the short arm of chromosome 6, which includes the SCA1 locus. A specific DNA gene probe will hybridize only with DNA derived from its target gene or gene fragment, and the resultant complex can be isolated and identified by techniques
5 known in the art.

In general, for detecting the presence of a DNA sequence located within the SCA1 gene, the genomic DNA is digested with a restriction endonuclease to obtain DNA fragments. The source of genomic DNA to be tested can be any biological specimen that contains DNA. Examples include specimen of blood,
10 semen, vaginal swabs, tissue, hair, and body fluids. The restriction endonuclease can be any that will cut the genomic DNA into fragments of double-stranded DNA having a particular nucleotide sequence. The specificities of numerous endonucleases are well known and can be found in a variety of publications, e.g. Maniatis et al.; Molecular Cloning: A Laboratory Manual; Cold Spring Harbor
15 Laboratory: New York (1982). That manual is incorporated herein by reference in its entirety. Preferred restriction endonuclease enzymes include *EcoRI*, *TaqI*, and *BstNI*. *EcoRI* is particularly preferred.

Diagnosis of the disease can alternatively involve the use of the polymerase chain reaction sequence amplification method (PCR) using novel
20 primers. U.S. Patent No. 4,683,195 (Mullis et al., issued July 28, 1987) describes a process for amplifying, detecting and/or cloning nucleic acid sequences. The method involves treating extracted DNA to form single-stranded complementary strands, treating the separate complementary strands of DNA with two oligonucleotide primers, extending the primers to form complementary extension
25 products that act as templates for synthesizing the desired nucleic acid molecule; and detecting the amplified molecule. More specifically, the method steps of treating the DNA with primers and extending the primers include the steps of: adding a pair of oligonucleotide primers, wherein one primer of the pair is substantially complementary to part of the sequence in the sense strand and the other
30 primer of each pair is substantially complementary to a different part of the same sequence in the complementary antisense strand; annealing the paired primers to the complementary molecule; simultaneously extending the annealed primers from a 3' terminus of each primer to synthesize an extension product complementary to the strands annealed to each primer wherein said extension products after separation

-15-

from the complement serve as templates for the synthesis of an extension product for the other primer of each pair; and separating said extension products from said templates to produce single-stranded molecules. Variations of the method are described in U.S. Patent No. 4,683,194 (Saiki et al., issued July 28, 1987). The polymerase chain reaction sequence amplification method is also described by Saiki
5 et al., *Science*, **230**, 1350-1354 (1985) and Scharf et al., *Science*, **324**, 163-166 (1986). The discussion of the these techniques in each of these references is incorporated herein by reference.

The primers are oligonucleotides, either synthetic or naturally
10 occurring, capable of acting as a point of initiating synthesis of a product complementary to the region of the DNA sequence containing the CAG repeating trinucleotides of the SCA1 locus of the short arm of chromosome 6. The primer includes a nucleotide sequence substantially complementary to a portion of a strand of an affected or a normal allele of a fragment (preferably a 3.36 kb *EcoRI*
15 fragment) of an SCA1 gene having a (CAG)_n region. The primer sequence has at least about 11 nucleotides, preferably at least about 16 nucleotides and no more than about 35 nucleotides. The primers are chosen such that they produce a primed product of about 70-350 base pairs, preferably about 100-300 base pairs. More preferably, the primers are chosen such that nucleotide sequence is substantially
20 complementary to a portion of a strand of an affected or a normal allele within about 150 nucleotides on either side of the (CAG)_n region, including directly adjacent to the (CAG)_n region.

Examples of preferred primers are shown by solid lines with arrowheads in Figure 3. The primers are thus selected from the group consisting of
25 CCGGAGCCCTGCTGAGGT (CAG-a), CCAGACGCCGGGACAC (CAG-b), AACTGGAAATGTGGACGTAC (Rep-1), CAACATGGGCAGTCTGAG (Rep-2), CCACCACTCCATCCCAGC (GCT-435), TGCTGGGCTGGTGGGGGG (GCT-214), CTCTCGGCTTTCTTGGTG (Pre-1), and GTACGTCCACATTTCCAGTT (Pre-2). These primers can be used in various combinations or with any other
30 primer that can be designed to hybridize to a portion of DNA of a fragment (preferably a 3.36 kb *EcoRI* fragment) of an SCA1 gene having a CAG repeat region. For example, the primer labeled Rep-2 can be combined with the primer labeled CAG-a, and the primer labeled CAG-b can be combined with the primer labeled Rep-1. More preferably the primers are the sets of primer pairs designed as

-16-

CAG-a/CAG-b, Rep-1/Rep-2, Rep-1/GCT-435, for example. These primer sets successfully amplify the CAG repeat units of interest using PCR technology. Alternatively, they can be used in various known techniques to sequence the SCA1 gene.

5 As stated previously, other methods of diagnosis can be used as well. They can be based on the isolation and identification of the repeat region of genomic DNA (CAG repeat region), cDNA (CAG repeat region), mRNA (GUC repeat region), and protein products (glutamine repeat region). These include, for example, using a variety of electrophoresis techniques to detect slight changes in the
10 nucleotide sequence of the SCA1 gene. Further nonlimiting examples include denaturing gradient electrophoresis, single strand conformational polymorphism gels, and nondenaturing gel electrophoresis techniques.

 The mapping and cloning of the SCA1 gene allows the definitive diagnosis of one type of the dominantly inherited ataxias using a simple blood test.
15 This represents the first step towards an unequivocal molecular classification of the dominant ataxias. A simple and reliable classification system for the ataxias is important because the clinical symptoms overlap extensively between the SCA1 and the non-SCA1 forms of the disease. Furthermore, a molecular test for the only known SCA1 mutation permits presymptomatic diagnosis of disease in known
20 SCA1 families and allows for the identification of sporadic or isolated CAG repeat expansions where there is no family history of the disease. Thus, the present invention can be used in family counseling, planning medical treatment, and in standard work-ups of patients with ataxia of unknown etiology.

25 B. Cloning

 Cloning of SCA1 DNA into the appropriate replicable vectors allows expression of the gene product, ataxin-1, and makes the SCA1 gene available for further genetic engineering. Expression of ataxin-1 or portions thereof, is useful because these gene products can be used as antigens to produce antibodies, as
30 described in more detail below.

1. Isolation of DNA

 DNA containing the SCA1 gene may be obtained from any cDNA library prepared from tissue believed to possess the SCA1 mRNA and to express it

-17-

at a detectable level. Preferably, the cDNA library is from human fetal brain or adult cerebellum. Optionally, the SCA1 gene may be obtained from a genomic DNA library or by *in vitro* oligonucleotide synthesis from the complete nucleotide or amino acid sequence.

5 Libraries are screened with appropriate probes designed to identify the gene of interest or the protein encoded by it. Preferably, for cDNA libraries, suitable probes include oligonucleotides that consist of known or suspected portions of the SCA1 cDNA from the same or different species; and/or complementary or homologous cDNAs or fragments thereof that consist of the same or a similar gene.
10 Optionally, for cDNA *expression* libraries (which express the protein), suitable probes include monoclonal or polyclonal antibodies that recognize and specifically bind to the SCA1 gene product, ataxin-1. Appropriate probes for screening *genomic* DNA libraries include, but are not limited to, oligonucleotides, cDNAs, or fragments thereof that consist of the same or a similar gene, and/or homologous
15 genomic DNAs or fragments thereof. Screening the cDNA or genomic library with the selected probe may be accomplished using standard procedures.

Screening cDNA libraries using synthetic oligonucleotides as probes is a preferred method of practicing this invention. The oligonucleotide sequences selected as probes should be of sufficient length and sufficiently unambiguous to
20 minimize false positives. The actual nucleotide sequence(s) of the probe(s) is usually designed based on regions of the SCA1 gene that have the least codon redundancy. The oligonucleotides may be degenerate at one or more positions, i.e., two or more different nucleotides may be incorporated into an oligonucleotide at a given position, resulting in multiple synthetic oligonucleotides. The use of
25 degenerate oligonucleotides is of particular importance where a library is screened from a species in which preferential codon usage is not known.

The oligonucleotide can be labeled such that it can be detected upon hybridization to DNA in the library being screened. A preferred method of labeling is to use ATP and polynucleotide kinase to radiolabel the 5' end of the
30 oligonucleotide. However, other methods may be used to label the oligonucleotide, including, but not limited to, biotinylation or enzyme labeling.

Of particular interest is the SCA1 nucleic acid that encodes a full-length mRNA transcript, including the complete coding region for the gene product,

-18-

ataxin-1. Nucleic acid containing the complete coding region can be obtained by screening selected cDNA libraries using the deduced amino acid sequence.

An alternative means to isolate the SCA1 gene is to use PCR methodology. This method requires the use of oligonucleotide primer probes that will hybridize to the SCA1 gene. Strategies for selection of PCR primer oligonucleotides are described below.

2. Insertion of DNA into Vector

The nucleic acid (e.g., cDNA or genomic DNA) containing the SCA1 gene is preferably inserted into a replicable vector for further cloning (amplification of the DNA) or for expression of the gene product, ataxin-1. Many vectors are available, and selection of the appropriate vector will depend on: 1) whether it is to be used for DNA amplification or for DNA expression; 2) the size of the nucleic acid to be inserted into the vector; and 3) the host cell to be transformed with the vector. Most expression vectors are "shuttle" vectors, i.e., they are capable of replication in at least one class of organism but can be transfected into another organism for expression. For example, a vector is cloned in *E. coli* and then the same vector is transfected into yeast or mammalian cells for expression even though it is not capable of replicating independently of the host cell chromosome. Each replicable vector contains various structural components depending on its function (amplification of DNA or expression of DNA) and the host cell with which it is compatible. These components are described in detail below.

Construction of suitable vectors employs standard ligation techniques known in the art. Isolated plasmids or DNA fragments are cleaved, tailored, and relegated in the form desired to generate the plasmids required. Typically, the ligation mixtures are used to transform *E. coli* K12 strain 294 (ATCC 31,446) and successful transformants are selected by ampicillin or tetracycline resistance where appropriate. Plasmids from the transformants are prepared, analyzed by restriction endonuclease digestion, and/or sequenced by methods known in the art. See, e.g., Messing et al., Nucl. Acids Res., 9, 309 (1981) and Maxam et al., Methods in Enzymology, 65, 499 (1980).

Optionally, DNA may also be amplified by direct insertion into the host genome. This is readily accomplished using *Bacillus* species as hosts, for example, by including in the vector a DNA sequence that is complementary to a

sequence found in *Bacillus* genomic DNA. Transfection of *Bacillus* with this vector results in homologous recombination with the genome and insertion of SCA1 DNA. However, the recovery of genomic DNA containing the SCA1 gene is more complex than that of an exogenously replicated vector because restriction enzyme
5 digestion is required to excise the SCA1 DNA.

Replicable cloning and expression vector components generally include, but are not limited to, one or more of the following: a signal sequence, an origin of replication, one or more marker genes, an enhancer element, a promoter and a transcription termination sequence.

10 *Vector component: signal sequence.* A signal sequence may be used to facilitate extracellular transport of a cloned protein. To this end, the SCA1 gene product, ataxin-1, may be expressed not only directly, but also as a fusion product with a heterologous polypeptide, preferably a signal sequence or other polypeptide having a specific cleavage site at the N-terminus of the cloned protein or
15 polypeptide. The signal sequence may be a component of the vector, or it may be a part of the SCA1 DNA that is inserted into the vector. The heterologous signal sequence selected should be one that is recognized and processed (i.e., cleaved by a signal peptidase) by the host cell. For prokaryotic host cells, a prokaryotic signal sequence may be selected, for example, from the group of the alkaline phosphatase, penicillinase, lpp or heat-stable intertoxin II leaders. For yeast secretion the signal
20 sequence used may be, for example, the yeast invertase, alpha factor, or acid phosphatase leaders. In mammalian cell expression, a native signal sequence may be satisfactory, although other mammalian signal sequences may be suitable, such as signal sequences from secreted polypeptides of the same or related species, as
25 well as viral secretory leaders, for example, the herpes simplex gD signal.

Vector component: origin of replication. Both expression and cloning vectors contain a nucleic acid sequence that enables the vector to replicate in one or more selected host cells. Generally, in cloning vectors this sequence is one that enables the vector to replicate independently of the host chromosomal DNA,
30 and includes origins of replication or autonomously replicating sequences. Such sequences are well known for a variety of bacteria, yeast and viruses. The origin of replication from the plasmid pBR322 is suitable for most Gram-negative bacteria, the 2m plasmid origin is suitable for yeast, and various viral origins (SV40, polyoma, adenovirus, VSV or BPV) are useful for cloning vectors in mammalian

-20-

cells. Generally, the origin of replication component is not needed for mammalian expression vectors (the SV40 origin may typically be used only because it contains the early promoter).

Vector component: marker gene. Expression and cloning vectors
5 may contain a marker gene, also termed a selection gene or selectable marker. This gene encodes a protein necessary for the survival or growth of transformed host cells grown in a selective culture medium. Host cells not transformed with the vector containing the selection gene will not survive in the culture medium. Typical selection genes encode proteins that: (a) confer resistance to antibiotics or other
10 toxins, e.g., ampicillin, neomycin, methotrexate, streptomycin or tetracycline; (b) complement auxotrophic deficiencies; or (c) supply critical nutrients not available from complex media, e.g., the gene encoding D-alanine racemase for *Bacilli*. One example of a selection scheme utilizes a drug to arrest growth of a host cell. Those cells that are successfully transformed with a heterologous gene express a protein
15 conferring drug resistance and thus survive the selection regimen.

An example of suitable selectable markers for mammalian cells are those that enable the identification of cells competent to take up the SCA1 nucleic acid, such as dihydrofolate reductase (DHFR) or thymidine kinase. The mammalian cell transformants are placed under selection pressure that only transformants are
20 uniquely adapted to survive by virtue of having taken up the marker. For example, cells transformed with the DHFR selection gene are first identified by culturing all the transformants in a culture medium that contains methotrexate, a competitive antagonist for DHFR. An appropriate host cell when wild-type DHFR is employed is the Chinese hamster ovary (CHO) cell line deficient in DHFR activity, prepared
25 and propagated as described by Urlaub et al., Proc. Natl. Acad. Sci. USA, 77, 4216 (1980). The transformed cells are then exposed to increased levels of methotrexate. This leads to the synthesis of multiple copies of the DHFR gene, and, concomitantly, multiple copies of the other DNA comprising the expression vectors, such as the SCA1 gene. This amplification technique can be used with any
30 otherwise suitable host, e.g., ATCC No. CCL61 CHO-K1, notwithstanding the presence of endogenous DHFR if, for example, a mutant DHFR gene that is highly resistant to methotrexate is employed. Alternatively, host cells (particularly wild-type hosts that contain endogenous DHFR) transformed or co-transformed with SCA1 DNA, wild-type DHFR protein, and another selectable marker such as

aminoglycoside 3' phosphotransferase (APH) can be selected by cell growth in a medium containing a selection agent for the selectable marker such as an aminoglycosidic antibiotic, e.g., kanamycin or neomycin. A suitable selection gene for use in yeast is the *trp1* gene present in the yeast plasmid YRp7 (Stinchcomb et al., Nature, 282, 39 (1979); Kingsman et al., Gene, 7, 141 (1979); or Tschemper et al., Gene, 10, 157 (1980)). The *trp1* gene provides a selection marker for a mutant strain of yeast lacking the ability to grow in tryptophan, for example, ATCC NO. 44076 or PEP4-1 (Jones, Genetics, 85, 12 (1977)). The presence of the *trp1* lesion in the yeast host cell genome then provides an effective environment for detecting transformation by growth in the absence of tryptophan. Similarly, *Leu2* deficient yeast strains (ATCC 20,622 or 38,626) are complemented by known plasmids bearing the *Leu2* gene.

Vector component: promoter. Expression and cloning vectors usually contain a promoter that is recognized by the host organism and is operably linked to the SCA1 nucleic acid. Promoters are untranslated sequences located upstream (5') to the start codon of a structural gene (generally within about 100 to 1000 bp) that control the transcription and translation of a particular nucleic acid sequence, such as the ataxin-1 nucleic acid sequence, to which they are operably linked. Such promoters typically fall into two classes, inducible and constitutive. Inducible promoters are promoters that initiate increased levels of transcription from DNA under their control in response to some change in culture conditions, e.g., the presence or absence of a nutrient or a change in temperature. In contrast, constitutive promoters produce a constant level of transcription of the cloned DNA segment.

At this time a large number of promoters recognized by a variety of potential host cells are well known in the art. Promoters are removed from their source DNA using a restriction enzyme digestion and inserted into the cloning vector using standard molecular biology techniques. Both the native SCA1 promoter sequence and many heterologous promoters can be used to direct amplification and/or expression of the SCA1 DNA. Heterologous promoters are preferred, as they generally permit greater transcription and higher yields of expressed protein as compared to the native promoter. Well-known promoters suitable for use with prokaryotic hosts include the beta-lactamase and lactose promoter systems, alkaline phosphatase, a tryptophan (*trp*) promoter system, and

hybrid promoters such as the tac promoter. Such promoters can be ligated to SCA1 DNA using linkers or adapters to supply any required restriction sites. Promoters for use in bacterial systems may contain a Shine-Dalgarno sequence for RNA polymerase binding.

5 Promoter sequences are known for eukaryotes. Virtually all eukaryotic genes have an AT-rich region located approximately 25 to 30 bp upstream from the site where transcription is initiated. Another sequence found 70 to 80 bases upstream from the start of transcription of many genes is the CXCAAT region where X may be any nucleotide. At the 3' end of most eukaryotic genes is an
10 AATAAA sequence that may be a signal for addition of the poly A tail to the 3' end of the coding sequence. All these sequences are suitably inserted into eukaryotic expression vectors. Examples of suitable promoting sequences for use with yeast hosts include the promoters for 3-phosphoglycerate kinase or other glycolytic enzymes, such as enolase, glyceraldehyde-3-phosphate dehydrogenase, hexokinase,
15 pyruvate decarboxylase, phosphofructokinase, glucose-6-phosphate isomerase, 3-phosphoglycerate mutase, pyruvate kinase, triosephosphate isomerase, phosphoglucose isomerase and glucokinase. Other yeast promoters, which are inducible promoters having the additional advantage of transcription controlled by growth conditions, are the promoter regions for alcohol dehydrogenase 2,
20 isocytochrome C, acid phosphatase, degradative enzymes associated with nitrogen metabolism, metallothionein, glyceraldehyde-3-phosphate dehydrogenase, and enzymes responsible for maltose and galactose utilization.

SCA1 transcription from vectors in mammalian host cells can be controlled, for example, by promoters obtained from the genomes of viruses such as
25 polyoma virus, fowlpox virus, adenovirus (such as Adenovirus 2), bovine papilloma virus, avian sarcoma virus, cytomegalovirus, a retrovirus, Hepatitis-B virus and most preferably Simian Virus 40 (SV40) (Fiers et al., Nature, 273, 113 (1978); Mulligan et al., Science, 209, 1422-1427 (1980); Pavlakis et al., Proc. Natl. Acad. Sci. USA, 78, 7398-7402 (1981)). Heterologous mammalian promoters (e.g., the
30 actin promoter or an immunoglobulin promoter) and heat-shock promoters can also be used, as can the promoter normally associated with the SCA1 sequence itself, provided such promoters are compatible with the host cell systems.

Vector component: enhancer element. Transcription of SCA1 DNA by higher eukaryotes can be increased by inserting an enhancer sequence into the

-23-

vector. Enhancers are *cis*-acting elements of DNA, usually having about 10 to 300 bp, that act on a promoter to increase its transcription. Enhancers are relatively orientation- and position-independent, having been found 5' and 3' to the transcription unit, within an intron as well as within the coding sequence itself.

5 Many enhancer sequences are now known from mammalian genes (globin, elastase, albumin, alpha-fetoprotein, and insulin). Typically, however, an enhancer from a eukaryotic cell virus will be used. Examples include the SV40 enhancer on the late side of the replication origin, the cytomegalovirus early promoter enhancer, the polyoma enhancer on the late side of the replication origin, and adenovirus
10 enhancers. The enhancer may be spliced into the vector at a position 5' or 3' to the SCA1 gene, but is preferably located at a site 5' of the promoter.

Vector component: transcription termination. Expression vectors used in eukaryotic host cells (yeast, fungi, insect, plant, animal, human or nucleated cells from other multicellular organisms) can also contain sequences necessary for
15 the termination of transcription and for stabilizing the mRNA. Such sequences are commonly available from the 5' and, occasionally, 3' untranslated regions of eukaryotic or viral DNAs or cDNAs. These regions can contain nucleotide segments transcribed as polyadenylated fragments in the untranslated portion of mRNA encoding ataxin-1.

20 Preferably, the pMALTM-2 vectors (New England Biolabs, Beverly, MA) are used to create the expression vector. These vectors provide a convenient method for expressing and purifying ataxin-1 produced from the cloned SCA1 gene. The SCA1 gene is inserted downstream from the *malE* gene of *E. coli*, which encodes maltose-binding protein (MBP) resulting in the expression of an MBP
25 fusion protein. The method uses the strong "tac" promoter and the *malE* translation initiation signals to give high-level expression of the cloned sequences, and a one-step purification of the fusion protein using MBP's affinity for maltose. The vectors express the *malE* gene (with or without its signal sequence) fused to the *lacZα* gene. Restriction sites between *malE* and *lacZα* are available for inserting the coding
30 sequence of interest. Insertion inactivates the β-galactosidase α-fragment activity of the *malE-lacZα* fusion, which results in a blue to white color change on Xgal plates when the construction is transformed into an α-complementing host such as TB1 (T.C. Johnston et al., *J. Biol. Chem.*, **261**, 4805-4811 (1986)) or JM107 (C. Yanisch-Perron et al., *Gene*, **33**, 103-119 (1985)). When present, the signal peptide on pre-

-24-

MBP directs fusion proteins to the periplasm. For fusion proteins that can be successfully exported, this allows folding and disulfide bond formation to take place in the periplasm of *E. coli*, as well as allowing purification of the protein from the periplasm. The vectors carry the *lac*^q gene, which codes for the Lac repressor protein. This keeps expression from P_{lac} low in the absence of isopropyl β-D-thiogalactopyranoside (IPTG) induction. The pMAL™-2 vectors also contain the sequence coding for the recognition site of the specific protease factor Xa, located just 5' to the polylinker insertion sites. This allows MBP to be cleaved from ataxin-1 after purification. Factor Xa cleaves after its four amino acid recognition sequence, so that few or no vector derived residues are attached to the protein of interest, depending on the site used for cloning.

Also useful are expression vectors that provide for transient expression in mammalian cells of SCA1 DNA. In general, transient expression involves the use of an expression vector that is able to replicate efficiently in a host cell, such that the host cell accumulates many copies of the expression vector and, in turn, synthesizes high levels of a desired polypeptide encoded by the expression vector. Transient expression systems, comprising a suitable expression vector and a host cell, allow for the convenient positive identification of polypeptides encoded by cloned DNAs, as well as for the rapid screening of such polypeptides for desired biological or physiological properties. Thus, transient expression systems are particularly useful in the invention for purposes of identifying analogs and variants of ataxin-1 that have wild-type or variant biological activity.

3. Host Cells

Suitable host cells for cloning or expressing the vectors herein are the prokaryote, yeast, or higher eukaryotic cells described above. Suitable prokaryotes include eubacteria, such as Gram-negative or Gram-positive organisms, for example, *E. coli*, *Bacilli* such as *B. subtilis*, *Pseudomonas* species such as *P. aeruginosa*, *Salmonella typhimurium*, or *Serratia marcescans*. One preferred *E. coli* cloning host is *E. coli* 294 (ATCC 31,446), although other strains such as *E. coli* B, *E. coli* X1776 (ATCC 31,537), and *E. coli* W3110 (ATCC 27,325) are suitable. These examples are illustrative rather than limiting. Preferably the host cell should secrete minimal amounts of proteolytic enzymes. Alternatively, *in vitro* methods of cloning, e.g., PCR or other nucleic acid polymerase reactions, are suitable.

-25-

In addition to prokaryotes, eukaryotic microbes such as filamentous fungi or yeast are suitable hosts for SCA1-encoding vectors. *Saccaromyces cerevisiae*, or common baker's yeast, is the most commonly used among lower eukaryotic host microorganisms. However, a number of other genera, species, and strains are commonly available and useful herein, such as *Schizosaccaromyces pombe*, *Kluyveromyces* hosts such as, e.g., *K. lactis*, *K. fragilis*, *K. bulgaricus*, *K. thermotolerans*, and *K. marxianus*, *Yarrowia*, *Pichia pastoris*, *Candida*, *Trichoderma reesia*, *Neurospora crassa*, and filamentous fungi such as, e.g., *Neurospora*, *Penicillium*, *Tolypocladium*, and *Aspergillus* hosts such as *A. nidulans*.

Suitable host cells for the expression of glycosylated ataxin-1 are derived from multicellular organisms. Such host cells are capable of complex processing and glycosylation activities. In principle, any higher eukaryotic cell culture is workable, whether from vertebrate or invertebrate culture. Examples of invertebrate cells include plant and insect cells. Numerous baculoviral strains and variants and corresponding permissive insect host cells from hosts such as *Spodoptera frugiperda* (caterpillar), *Aedes aegypti* (mosquito), *Aedes albopictus* (mosquito), *Drosophila melanogaster* (fruitfly), and *Bombyx mori* have been identified. See, e.g., Luckow et al., Bio/Technology, 6, 47-55 (1988); Miller et al., Genetic Engineering, 8, 277-279 (1986); and Maeda et al., Nature, 315, 592-594 (1985). A variety of viral strains for transfection are publicly available, e.g., the L-1 variant of *Autographa californica* NPV and the Bm-5 strain of *Bombyx mori* NPV, and such viruses may be used as the virus herein according to the present invention, particularly for transfection of *Spodoptera frugiperda* cells.

Plant cell cultures of cotton, corn, potato, soybean, petunia, tomato, and tobacco can be utilized as hosts. Typically, plant cells are transfected by incubation with certain strains of the bacterium *Agrobacterium tumefaciens*, which has been previously manipulated to contain the SCA1 DNA. During incubation of the plant cell culture with *A. tumefaciens*, the SCA1 DNA is transferred to the plant cell host such that it is transfected, and will, under appropriate conditions, express the SCA1 DNA. In addition, regulatory and signal sequences compatible with plant cells are available, such as the nopaline synthase promoter and polyadenylation signal sequences. Depicker et al., J. Mol. Appl. Gen., 1, 561 (1982).

Vertebrate cells can also be used as hosts. Propagation of vertebrate cells in culture (tissue culture) has become a routine procedure in recent years.

-26-

Examples of useful mammalian host cell lines are monkey kidney CV1 line transformed by SV40 (CAS-7, ATCC CRL 1651); human embryonic kidney line (293 or 293 cells subcloned for growth in suspension culture, Graham et al., J. Gen. Virol., 36, 59 (1977)); baby hamster kidney cells (BHK, ATCC CCL 10); Chinese hamster ovary cells/-DHFR (CHO, Urlaub and Chasin, Proc. Natl. Acad. Sci. USA, 77, 4216 (1980)); mouse sertoli cells (TM4, Mather, Biol. Reprod., 23, 243-251 (1980)); monkey kidney cells (CV1 ATCC CCL 70); African green monkey kidney cells (VERO-76, ATCC CRL-1587); human cervical carcinoma cells (HELA, ATCC CCL 2); canine kidney cells (MDCK, ATCC CCL 34); buffalo rat liver cells (BRL 3A, ATCC CRL 1442); human lung cells (W138, ATCC CCL 75); human liver cells (Hep G2, HB 8065); mouse mammary tumor (MMT 060562, ATCC CCL 51); TRI cells (Mather et al., Annals N.Y. Acad. Sci., 383, 44-68 (1982)); MRC 5 cells; FS4 cells; and a human hepatoma line (Hep G2).

4. Transfection and transformation

Host cells are transfected and preferably transformed with the above-described expression or cloning vectors of this invention and cultured in conventional nutrient media modified as appropriate for inducing promoters, selecting transformants, or amplifying the genes encoding the desired sequences.

Transfection refers to the taking up of an expression vector by a host cell whether or not any coding sequence are in fact expressed. Numerous methods of transfection are known to the ordinarily skilled artisan, for example, the calcium phosphate precipitation method and electroporation are commonly used. Successful transfection is generally recognized when any indication of the operation of the vector occurs within the host cell.

Transformation means introducing DNA into an organism so that the DNA is replicable, either as an extrachromosomal element or by chromosomal integrant. Depending on the host cell used, transformation is done using standard techniques appropriate to such cells. Calcium chloride is generally used for prokaryotes or other cells that contain substantial cell-wall barriers. Infection with *Agrobacterium tumefaciens* can be used for transformation of certain plant cells. For mammalian cells without cell walls, the calcium phosphate precipitation method of Graham et al., Virology, 52, 456-457 (1978) is preferred. Transformations into yeast are typically carried out according to the method of Van Solingen et al., J.

-27-

Bact., 130, 946 (1977) and Hsiao et al., Proc. Natl. Acad. Sci. (USA), 78 3829 (1979). However, other methods for introducing DNA into cells such as by nuclear injection, electroporation, or protoplast fusion may also be used.

5 5. Cell Culture

Prokaryotic cells used to produce the SCA1 gene product, ataxin-1, are cultured in suitable media, as described generally in Sambrook et al. The mammalian host cells used to produce the SCA1 gene product may be cultured in a variety of media. Commercially available media such as Hams F10 (Sigma),
10 Minimal Essential Medium (MEM, Sigma), RPMI-1640 (Sigma), and Dulbecco's Modified Eagle's Medium (DMEM, Sigma) are suitable for culturing the host cells. These media may be supplemented as necessary with hormones and/or other growth factors (such as insulin, transferrin, or epidermal growth factor), salts (such as sodium chloride, calcium, magnesium, and phosphate), buffers (such as HEPES),
15 nucleosides (such as adenosine and thymidine), antibiotics (such as Gentamycin™ drug), trace elements (defined as inorganic compounds usually present at final concentrations in the micromolar range), and glucose or an equivalent energy source. Any other necessary supplements may also be included at appropriate concentrations that would be known to those skilled in the art. The culture
20 conditions, such as temperature, pH, and the like, are those previously used with the host cell selected for expression, and will be apparent to the ordinarily skilled artisan. The host cells referred to in this disclosure encompass in *in vitro* culture as well as cells that are within a host animal.

25 C. Protein

The SCA1 gene encodes a novel protein, ataxin-1, a representative example of which is shown in Figure 15 with an estimated molecular weight of about 87 kD. It is to be understood that ataxin-1 represents a set of proteins produced from the SCA1 gene with its unstable CAG region. Ataxin-1 can be
30 produced from cell cultures. With the aid of recombinant DNA techniques, synthetic DNA and cDNA coding for ataxin-1 can be introduced into microorganisms which can then be made to produce the peptide. It is also possible to manufacture ataxin-1 synthetically, in a manner such as is known for peptide syntheses.

-28-

Ataxin-1 is preferably recovered from the culture medium as a cytosolic polypeptide, although it can also be recovered as a secreted polypeptide when expressed with a secretory signal.

Ataxin-1 can be purified from recombinant cell proteins or polypeptides to obtain preparations that are substantially homogenous as ataxin-1. As a first step, the culture medium or lysate is centrifuged to remove particulate cell debris. The membrane and soluble protein fractions are then separated. The ataxin-1 may then be purified from the soluble protein fraction and from the membrane fraction of the culture lysate, depending on whether the ataxin-1 is membrane bound. If necessary, ataxin-1 is further purified from contaminant soluble proteins and polypeptides, with the following procedures being exemplary of suitable purification procedures: by fractionation on immunoaffinity or ion-exchange columns; ethanol precipitation; reverse phase HPLC; chromatography on silica or on a cation-exchange resin such as DEAE; chromatofocusing; SDS-PAGE; ammonium sulfate precipitation; gel filtration using, for example, Sephadex G-75; ligand affinity chromatography, using, e.g., protein A Sepharose columns to remove contaminants such as IgG.

Ataxin-1 variants in which residues have been deleted, inserted, or substituted are recovered in the same fashion as native ataxin-1, taking account of any substantial changes in properties occasioned by the variation. For example, preparation of a ataxin-1 fusion with another protein or polypeptide, e.g., a bacterial or viral antigen, facilitates purification; an immunoaffinity column containing antibody to the antigen can be used to adsorb the fusion polypeptide. Immunoaffinity columns such as a rabbit polyclonal ataxin-1 column can be employed to absorb the ataxin-1 variant by binding it to at least one remaining immune epitope. Alternatively, the ataxin-1 may be purified by affinity chromatography using a purified ataxin-1-IgG coupled to a (preferably) immobilized resin such as Affi-Gel 10 (Bio-Rad, Richmond, CA) or the like, by means well-known in the art. A protease inhibitor such as phenyl methyl sulfonyl fluoride (PMSF) also may be useful to inhibit proteolytic degradation during purification, and antibiotics may be included to prevent the growth of adventitious contaminants.

Covalent modifications of ataxin-1 are included within the scope of this invention. Both native ataxin-1 and amino acid sequence variants of the ataxin-1 may be covalently modified. Covalent modifications included within the scope of

-29-

this invention are those producing one or more ataxin-1 fragments. Ataxin-1 fragments having any number of amino acid residues may be conveniently prepared by chemical synthesis, by enzymatic or chemical cleavage of the full-length or variant ataxin-1 polypeptide, or by cloning and expressing only portions of the
5 SCA1 gene. Other types of covalent modifications of ataxin-1 or fragments thereof are introduced into the molecule by reacting targeted amino acid residues of the ataxin-1 or fragments thereof with a derivatizing agent capable of reacting with selected side chains or the N- or C-terminal residues.

For example, cysteinyl residues most commonly are reacted with α -
10 haloacetates (and corresponding amines), such as iodoacetic acid or iodoacetamide, to give carboxymethyl or carboxyamidomethyl derivatives. Cysteinyl residues also are derivatized by reaction with bromotrifluoroacetone, α -bromo- β -(5-imidozoyl)propionic acid, iodoacetyl phosphate, N-alkylmaleimides, 3-nitro-2-pyridyl disulfide, methyl 2-pyridyl disulfide, *p*-chloromercuribenzoate, 2-
15 chloromercuri-4-nitrophenol, or chloro-7-nitrobenzo-2-oxa-1,3-diazole.

Histidyl residues are derivatized by reaction with diethylpyrocarbonate *p*-bromophenacyl. Lysinyl and amino terminal residues are derivatized with succinic or other carboxylic acid anhydrides and imidoesters such as methyl picolinimide; pyridoxal phosphate; pyridoxal; chloroborohydride;
20 trinitrobenzenesulfonic acid; O-methylisourea; 2,4-pentanedione; and transaminase-catalyzed reaction with glyoxylate. Arginyl residues are modified by reaction with phenylglyoxal, 2,3-butanedione, 1,2-cyclohexanedione, and ninhydrin, among others.

Specific modification of tyrosyl residues may be made, with
25 particular interest in introducing spectral labels into tyrosyl residues by reaction with aromatic diazonium compounds or tetranitromethane. Most commonly, N-acetylimidazole and tetranitromethane are used to form O-acetyl tyrosyl species and 3-nitro derivatives, respectively. Tyrosyl residues are iodinated using ^{125}I or ^{131}I to prepared labeled proteins for use in radioimmunoassay, the chloramine T method
30 described above being suitable.

Carboxyl side groups (aspartyl or glutamyl) are selectively modified by reaction with carbodiimides ($\text{R-N}=\text{C}=\text{N-R}'$), where R and R' are different alkyl groups, such as 1-cyclohexyl-3-(2-morpholinyl-4-ethyl)carbodiimide or 1-ethyl-3-(4-azonia-4,4-dimethylpentyl)carbodiimide. Furthermore, aspartyl and glutamyl

-30-

residues are converted to asparaginyl and glutaminyl residues by reaction with ammonium ions.

Derivatization with bifunctional agents is useful for crosslinking ataxin-1 to a water-insoluble support matrix or surface for use in the method for
5 purifying anti-ataxin-1 antibodies, and vice versa. Commonly used crosslinking agents include, *e.g.*, 1,1-bis(diazoacetyl)-2-phenylethane, glutaraldehyde, and N-hydroxysuccinimide esters, for example, esters with 4-azidosalicylic acid, homobifunctional imidoesters, including disuccinimidyl esters such as 3,3'-dithiobis(succinimidylpropionate), and bifunctional maleimides such as bis-N-
10 maleimido-1,8-octane. Derivatizing agents such as methyl-3-[(*p*-azidophenyl)dithio]propionimide yield photoactivatable intermediates that are capable of forming crosslinks in the presence of light. Alternatively, reactive water-insoluble matrices such as cyanogen bromide-activated carbohydrates and the reactive substrates are employed for protein immobilization.

15 Glutaminyl and asparaginyl residues are frequently deamidated to the corresponding glutamyl and aspartyl residues, respectively. These residues are deamidated under neutral or basic conditions. The deamidated form of these residues falls within the scope of this invention.

Other modifications include hydroxylation of proline and lysine,
20 phosphorylation of hydroxyl groups of seryl or threonyl residues, methylation of the α -amino groups of lysine, arginine, and histidine side chains, acetylation of the N-terminal amine, amidation of any C-terminal carboxyl group, and glycosylation of any suitable residue.

25 **D. Antibodies**

The present invention also relates to polyclonal or monoclonal antibodies raised against ataxin-1 or ataxin-1 fragments (preferably fragments having 8-40 amino acids, more preferably 10-20 amino acids, that form the surface of the folded protein), or variants thereof, and to diagnostic methods based on the
30 use of such antibodies, including but not limited to Western blotting and ELISA (enzyme-linked immunosorbant assay).

Polyclonal antibodies to the SCA1 polypeptide generally are raised in animals by multiple subcutaneous (sc) or intraperitoneal (ip) injections of ataxin-1, ataxin-1 fragments, or variants thereof, and an adjuvant. The polypeptide can be a

-31-

cloned gene product or a synthetic molecule. Preferably, it corresponds to a position in the protein sequence that is on the surface of the folded protein and is thus likely to be antigenic. It may be useful to conjugate the SCA1 polypeptide (including fragments containing a specific amino acid sequence) to a protein that is immunogenic in the species to be immunized, *e.g.*, keyhole limpet hemocyanin, serum albumin, bovine thyroglobulin, or soybean trypsin inhibitor using a bifunctional or derivatizing agent, for example, maleimidobenzoyl sulfosuccinimide ester (conjugation through cysteine residues), N-hydroxysuccinimide (through lysine residues), glutaraldehyde, succinic anhydride, SOCl_2 , or $\text{R}^1\text{N}=\text{C}=\text{NR}$, where R and R^1 are different alkyl groups. Conjugates also can be made in recombinant cell culture as protein fusions. Also, aggregating agents such as alum are used to enhance the immune response.

The route and schedule of immunizing a host animal or removing and culturing antibody-producing cells are variable and are generally in keeping with established and conventional techniques for antibody stimulation and production. While mice are frequently employed as the host animal, it is contemplated that any mammalian subject including human subjects or antibody-producing cells obtained therefrom can be manipulated according to the processes of this invention to serve as the basis for production of mammalian, including human, hybrid cell lines. Preferably, rabbits are used to raise antibodies against ataxin-1.

Animals are typically immunized against the immunogenic conjugates or derivatives by combining about 10 μg to about 1 mg of ataxin-1 with about 2-3 volumes of Freund's complete adjuvant and injecting the solution intradermally at multiple sites. About one month later the animals are boosted with about 1/5 to about 1/10 the original amount of conjugate in Freund's complete adjuvant (or other suitable adjuvant) by subcutaneous injection at multiple sites. About 7 to 14 days later animals are bled and the serum is assayed for anti-ataxin-1 polypeptide titer.

Serum antibodies (IgG) are purified via protein purification protocols that are well known in the art. Antibody/antigen reactivity is analyzed using Western blotting, wherein suspected antigens are blotted to a nitrocellulose filter, exposed to potential antibodies and allowed to hybridize under defined conditions. See Gershoni et al., Anal. Biochem., 131, 1-15 (1983). The protein antigens can

-32-

then be sequenced using standard sequencing methods directly from the antibody/antigen complexes on the nitrocellulose support.

Monoclonal antibodies are prepared by recovering immune cells - typically spleen cells or lymphocytes from lymph node tissue - from immunized animals (usually mice) and immortalizing the cells in conventional fashion, e.g., by fusion with myeloma cells. The hybridoma technique described originally by Kohler et al., Eur. J. Immunol., 6, 511 (1976) has been widely applied to produce hybrid cell lines that secrete high levels of monoclonal antibodies against many specific antigens. It is possible to fuse cells of one species with another. However, it is preferable that the source of the immunized antibody-producing cells and the myeloma be from the same species. While mouse monoclonal antibodies are routinely used, the present invention is not so limited. In fact, although mouse monoclonal antibodies are typically used, human antibodies may be used and may prove to be preferable. Such antibodies can be obtained by using human hybridomas. Cote et al.; Monoclonal Antibodies and Cancer Therapy; A.R. Liss, Ed.; p. 77 (1985).

The secreted antibody is recovered from tissue culture supernatant by conventional methods such as precipitation, ion exchange chromatography, affinity chromatography, or the like. The antibodies described herein are also recovered from hybridoma cell cultures by conventional methods for purification of IgG or IgM, as the case may be, that heretofore have been used to purify these immunoglobulins from pooled plasma, e.g., ethanol or polyethylene glycol precipitation procedures. The purified antibodies are sterile filtered, and optionally are conjugated to a detectable marker such as an enzyme or spin label for use in diagnostic assays of the ataxin-1 in test samples.

Techniques for creating recombinant DNA versions of the antigen-binding regions of antibody molecules (known as Fab fragments), which bypass the generation of monoclonal antibodies, are encompassed within the practice of this invention. Antibody-specific messenger RNA molecules are extracted from immune system cells taken from an immunized animal, transcribed into complementary DNA (cDNA), and the cDNA is cloned into a bacterial expression system.

The anti-ataxin-1 antibody preparations of the present invention are specific to ataxin-1 and do not react immunochemically with other substances in a

-33-

manner that would interfere with a given use. For example, they can be used to screen for the presence of ataxin-1 in tissue extracts to determine tissue-specific expression levels of ataxin-1.

The present invention also encompasses an immunochemical assay that involves subjecting antibodies directed against ataxin-1 to reaction with the ataxin-1 present in a sample to thus form an (ataxin-1/anti-ataxin-1) immune complex, the formation and amount of which are measures - qualitative and quantitative, respectively - of the ataxin-1 presence in the sample. The addition of other reagents capable of biospecifically reacting with constituents of the protein/antibody complex, such as anti-antibodies provided with analytically detectable groups, facilitates detection and quantification of ataxin-1 in biological samples, and is especially useful for quantitating the level of ataxin-1 in biological samples. Ataxin-1/anti-ataxin-1 complexes can also be subjected to amino acid sequencing using methods well known in the art to determine the length of a polyglutamine region and thereby provide information about likelihood of affliction with spinocerebellar ataxia and likely age of onset. Competitive inhibition and non-competitive methods, precipitation methods, heterogeneous and homogeneous methods, various methods named according to the analytically detectable group employed, immunoelectrophoresis, particle agglutination, immunodiffusion and immunohistochemical methods employing labeled antibodies may all be used in connection with the immune assay described above.

The invention has been described with reference to various specific and preferred embodiments and will be further described by reference to the following detailed examples. It is understood, however, that there are many extensions, variations, and modifications on the basic theme of the present invention beyond that shown in the examples and detailed description, which are within the spirit and scope of the present invention.

Experimental Section

I. The Gene for SCA1 Maps Centromeric to D6S89

To confirm the position of SCA1 with respect to D6S89 and to identify closer flanking markers, two dinucleotide repeat polymorphisms D6S109 and D6S202 were used. Using YAC clones isolated in the D6S89 region, three additional dinucleotide repeat polymorphisms were identified, one of which (AM10GA) showed no recombination with SCA1 and confirmed that D6S89 is telomeric to SCA1. The dinucleotide repeat at D6S109 revealed six recombination events with SCA1 and determined D6S109 to be the other flanking marker at the centromeric end. Linkage analysis, physical mapping data as discussed below, and analysis of recombination events demonstrated that the order of markers is as follows: Centomere - D6S109 - AM10GA/SCA1 - D6S89 - SB1 - LR40 - D6S202 - Telomere.

A. Materials and Methods

1. SCA1 Kindreds

Nine large SCA1 families were used in the present study. Clinical findings and linkage data demonstrating that these families segregated SCA1 have been previously reported. See, J.F. Jackson et al., N. Engl. J. Med., **296**, 1138-1141 (1977); B.J.B. Keats et al., Am. J. Hum. Genet., **49**, 972-977 (1991); L.P.W. Ranum et al., Am. J. Hum. Genet., **49**, 31-41 (1991); and H.Y. Zoghbi et al., Am. J. Hum. Genet., **49**, 23-30 (1991). Analysis of polymorphisms at the loci D6S109, AM10GA, SB1, LR40, and D6S202 was performed on individuals from these kindreds.

The Houston (TX-SCA1) kindred included 106 individuals, of whom 57 (25 affected) were genotyped. See, H.Y. Zoghbi et al., Ann. Neurol., **23**, 580-584 (1988). Patients symptomatic at the time of exam, as well as asymptomatic individuals who had both a symptomatic child and a symptomatic parent, were classified as "affected." In this kindred, a deceased individual previously assigned as affected (from family history data) was reassigned an unknown status after review of medical records. This reassignment eliminated what was previously thought to be a recombination event between SCA1 and D6S89 in the TX-SCA1 kindred. To maximize the amount of information available for linkage analysis, the two chromosomes 6 in somatic cell hybrids for 15 affected individuals and one

unaffected individual from the TX-SCA1 kindred were separated. See, H.Y. Zoghbi et al., Am. J. Hum. Genet., 44, 255-263 (1989). The Louisiana (LA-SCA1) kindred included 50 individuals of whom 26 (8 affected) were genotyped. See, B.J.B. Keats et al., Am. J. Hum. Genet., 49, 972-977 (1991). The Minnesota (MN-SCA1) kindred included 175 individuals, of whom 106 (17 affected) were genotyped. See, J.L. Haines et al., Neurology, 34, 1542-1548 (1984); and L.P.W. Ranum et al., Am. J. Hum. Genet., 49, 31-41 (1991). The Michigan (MI-SCA1) kindred included 201 individuals, of whom 127 (25 affected) were genotyped. See, H.E. Nino et al., Neurology, 30, 12-20 (1980). The Mississippi (MS-SCA1) kindred included 84 individuals, of whom 37 (17 affected) were genotyped. See, J.F. Jackson et al., N. Engl. J. Med., 296, 1138-1141 (1977).

Four Italian families segregating SCA1 were analyzed; their clinical phenotype and HLA linkage data were reported previously. See, M. Spadaro et al., Acta Neurol. Scand., 85, 257-265 (1992). Three families originated in the Calabria Region (Southern Italy): family IT-P with 135 members of whom 80 (21 affected) were genotyped; for computational reasons, the family was subdivided into 3 different pedigrees (RM, VI, and FB) and only one of the 3 consanguinity loops was considered; family IT-NS, with 43 members of whom 27 (7 affected) were typed; family IT-NS with 51 members of whom 16 (3 affected) were typed. The fourth family, IT-MR, originated from Latium and consisted of 17 individuals of whom 10 (4 affected) were genotyped.

2. CEPH Families

The 40 CEPH reference families were genotyped at the D9S109, LR40 and D6S202 loci in order to provide a large number of informative meioses for marker-marker linkage analyses. Markers AM10GA and SB1 flank D6S89, having been isolated from a yeast artificial chromosome (YAC) contig built bidirectionally from D6S89 (see below). A subset of 18 CEPH families which defined 26 recombinants between D6S109 and D6S89 was genotyped at AM10GA and SB1 in order to determine the order of AM10GA, D6S89 and SB1 with respect to D6S109.

3. Cloning of Sequences Containing Dinucleotide Repeats

The identification and description of polymorphic dinucleotide repeats at the D6S109 and D6S202 loci have been previously reported. See, L.P.W. Ranum

-36-

et al., Nucleic Acids Res., **19**, 1171 (1991); and F. LeBorgne-Demarquoy et al., Nucleic Acids Res., **19**, 6060 (1991).

DNA fragments containing dinucleotide repeats were cloned at LR40 and SB1 from yeast artificial chromosome (YAC) clones at the LR40 and FLB1 loci, respectively (see below). DNA from each YAC clone was amplified in a 50 µl reaction containing 20 ng DNA, a single Alu primer (see below), 50 mM KCl, 10 mM Tris-Cl pH 8.3, 1.25 mM MgCl₂, 200 or 250 µM dNTPs, 0.01% (w/v) gelatin, and 1.25 units *Thermus aquaticus* DNA polymerase (Taq polymerase--Perkin Elmer, Norwalk, CT). For amplification of FLB1 YAC DNA, a primer complementary to the 5' end of the Alu consensus sequence (Oncor Laboratories, Gaithersburg, MD), designated SAL1, was used = 5'-AGGAGTGAGCCACCGCACCCAGCC-3' at a final concentration of 0.6 µM. For amplification of LR40 YAC DNA, 0.2 µM primer PDJ34 was used. See, C. Breukel et al., Nucleic Acids Res., **18**, 3097 (1990). Samples were overlaid with mineral oil, denatured at 94°C for 5 minutes, then subjected to 30 cycles of 1 minute 94°C denaturation, 1 minute 55°C annealing, and 5 minutes 72°C extension. The last extension step was lengthened to 10 minutes. Electrophoresis of 15 µl of PCR products was performed on a 1.5% agarose gel, which was Southern blotted and hybridized with a probe prepared by random-hexamer-primed labelling of synthetic poly(dG-dT)-poly(dA-dC) (Pharmacia, Piscataway, NJ) using [α -³²P]dCTP, as described by A.P. Feinberg et al., Anal. Biochem., **137**, 266-267 (1984). Fragments hybridizing with the dinucleotide repeat probe were identified and were subsequently purified by electrophoresis on a low-melt agarose gel. Fragments were excised and reamplified by PCR as above.

For LR40, reamplified DNA was repurified by low-melt gel electrophoresis, and DNA extracted from excised bands by passage through a glasswool spin column as described by D.M. Heery et al., Trends Genet., **6**, 173 (1990). A purified 1.2-kb fragment was cloned into pBluescript plasmid modified as a "T-vector" as described by D. Marchuck et al., Nucleic Acids Res., **19**, 1154 (1990). From this clone, a 0.6-kb *Hinc*II restriction fragment containing a GT repeat was subcloned into pBluescript plasmid, and sequenced on an Applied Biosystems, Inc. (Foster City, CA) automated sequencer.

-37-

For SB1, a reamplified 1-kb fragment was ethanol precipitated and blunt-end cloned into pBluescript plasmid. Plasmid DNA was isolated and PCR amplified in one reaction with M13 Reverse primer plus BamGT primer (5'-CCCGGATCCTGTGTGTGTGTGTGTGTG-3') and in a second reaction M13
5 Universal primer and BamCA primer (5'-CCCGGATCCACACACACACACACAC-3'). See, C.A. Feener et al., Am. J. Hum. Genet., 48, 621-627 (1991). PCR conditions were as above except primers were used at 1 μ M concentration; 2.5 units Taq polymerase and approximately 30 ng DNA were used per reaction, with final reaction volumes of 100 μ l, and an
10 annealing temperature of 50°C. Products were precipitated, resuspended, and digested with *Bam*H1 (product of Universal primer reaction) or *Bam*H1 and *Hinc*II (product of Reverse primer reaction). These two fragments were cloned into pBluescript plasmid and sequenced as above.

Dinucleotide repeats were cloned at AM10 from a YAC containing
15 this locus. A λ FixII library was constructed using DNA from this yeast clone, and human clones were identified by filter hybridization using human placental DNA as a probe. A gridded array of these human clones was grown, and filters containing DNA from these clones were hybridized with a 32 P-labelled poly(dG-dT)-poly(dA-dC3) probe as described above. DNA was prepared from positive clones, digested
20 with various restriction enzymes, and analyzed by agarose gel electrophoresis. Southern blotting and hybridization were carried out with the poly(dG-dT)-poly(dA-dC) probe. A 1-kb fragment hybridizing with the dinucleotide repeat probe was identified, clones into M13, and sequenced.

25 4. PCR Analysis

Primer sequences and concentrations, and PCR cycle times used for amplification of dinucleotide repeat sequences from human genomic DNA are presented in Table 1. For the LR40 polymorphism, primer set "A" was used for analysis of the TX-SCA1, LA-SCA1, and MS-SCA1 kindreds, while primer set "B"
30 was used for all other kindreds. Buffer compositions were as follows: 50 mM KCl, 10 mM Tris-Cl pH 8.3, 1.25 mM $MgCl_2$ (1.5 mM $MgCl_2$ for AM10GA), 250 μ M dNTPs (200 μ M dNTPs for AM10GA), 0.01% (w/v) gelatin, and 0.5 - 0.625 unit Taq polymerase. For the LR40 analysis, 2% formamide was included in the PCR buffer. When primer set B was used for LR40 analysis, 125 μ M dNTPs, 1.5 mM

-38-

MgCl₂, and 1 unit Taq polymerase were used. All reaction volumes were 25 µl and contained 40 ng genomic DNA. Four microliters of each reaction was mixed with 2 µl formamide loading buffer, denatured at 90-100°C for 3 minutes, cooled on ice, and 2-4 µl was used for electrophoresis on a 4% or 6% polyacrylamide/7.65 M urea sequencing gel for 2-3 hours at 1100 V. PCR assay conditions have been reported previously for D6S202 and D6S109. See, L.P.W. Ranum et al., Nucleic Acids Res., 19, 1171 (1991); and F. LeBorgne-Demarquoy et al., Nucleic Acids Res., 19, 6060 (1991).

-39-

Table 1.
Primers and PCR conditions for amplification of
dinucleotide repeat sequences

<u>Marker/Type</u>	<u>Primers^a</u>	<u>PCR</u>	
		<u>Steps</u>	<u>Cycles</u>
AM10GA/(GA) _n	AAGTCAGCCTCTACTCTTTGT	94°C for 30 sec.	
	TGA		
	CTTGGAGCAGTCTGTAGGGAG	55°C for 30 sec. 72°C for 30 sec.	30
SB1/(GT) _n	TGAAGTGATGTGCTCTGTTC	94°C for 60 sec.	
	AAAGGGGTAGAGGAAATGAG	60°C for 60 sec. 72°C for 60 sec.	30
LR40/(GT) _n set A	AGGAGAGGGGTCATGAGTTG	94°C for 60 sec.	
	GGCTCATGAATACATTACATG		
	AAG	58°C for 60 sec. 72°C for 60 sec.	25
LR40/(GT) _n set B	CTCATTACCTTAGAGACAAA		
	TGGATAG	94°C for 60 sec.	
	ATGGTATAGGGATTTNCCAA	60°C for 60 sec. 72°C for 45 sec.	27
	ACCTG		

^aPrimers are shown as 5' to 3' sequence. The first primer of each pair was end-labelled with γ -³²P ATP and polynucleotide kinase. Primer concentrations were 1 mM.

5. SCA1 Linkage Analysis

The D6S109, AM10GA, D6S89, SB1, LR40 and D6S202 markers were analyzed for linkage to SCA1 using the computer program LINKAGE version 5.1 which includes the MLINK, ILINK, LINKMAP, CLODScore and CMAP programs. See, G.M. Lathrop et al., Proc. Natl. Acad. Sci. USA, **81**, 3443-3446 (1984). Age dependent penetrance classes were assigned independently for each of the families included in the analysis. Marker alleles were recoded to reduce the number of alleles segregating in a family to four, five or six alleles to simplify the analysis. The allele frequencies for the various markers were based on the frequencies of the alleles among the spouses in each family and were determined separately for the two American black kindreds, for the Italian kindreds, and for the Caucasian kindreds from Minnesota, Michigan, and Mississippi, with the following exception - the allele frequencies for D6S109 in the MI and MN kindreds were based on the frequencies of the alleles in the CEPH families.

Maximum LOD scores for the various markers were calculated with the MLINK program by running each of the analyses separately for the various families, at theta values with increments of 0.0005 to 0.001, and then adding the values of each of the kindreds. The analyses were done separately to ensure that the allele frequencies for the various markers were representative for each of the ethnically diverse families. As a control, the recombination fractions at the maximum lod scores (Z_{\max}) between each marker and SCA1 were calculated using the ILINK program after the allele frequencies for each marker were set equal to one another. In all cases the recombination frequencies were the same and Z_{\max} values were very similar to those reported in Table 5 below.

6. CEPH Linkage Analysis

Forty CEPH families were typed for the GT repeat markers D6S109, D6S202 and LR40. The original alleles were recoded to five alleles. The SB1 and AM10 markers were typed in a subset of the CEPH panel which defined 26 recombinants from 18 different families between D6S109 and D6S89. The CLODScore program was used for the two-point analyses and CMAP was used for the three-and four-point analyses. For the three-point and four-point analyses, the interval between the mapped markers was fixed based on the two point $\theta_m = \theta_f$ results. The likelihood of the location of the test locus (SCA1) was calculated at 10

different positions within each interval. The test for sex difference in the Θ values was performed using a χ^2 statistic, with $\chi^2 = 2(\ln 10)[Z(\theta_m, \theta_f) - Z(\theta = \theta_m = \theta_f)]$, where $Z(\theta_m, \theta_f)$ is the overall Z_{\max} for arbitrary θ_m and θ_f , while $Z(\theta = \theta_m = \theta_f)$ is the Z_{\max} constrained to $\theta_m = \theta_f$. Under homogeneity (H1), χ^2 approximates a χ^2 with 1 d.f. Rejection of homogeneity occurs when $\chi^2 > 3.84$.

B. Results

1. Dinucleotide Repeat Cloning and Sequencing and Analysis

Dinucleotide repeats SB1 and LR40 were amplified directly from YAC clones by *Alu*-primed PCR and the dinucleotide repeat containing fragments were identified by hybridization. The PCR products were cloned either directly or by further amplification using tailed poly(GT) or poly(CA) primers paired with an *Alu* primer. In addition, two dinucleotide repeats were subcloned from a lambda phage clone from a library constructed from a YAC at the AM10 locus.

Dinucleotide repeats from the SB1, LR40, and AM10 loci were sequenced. At LR40, the cloned repeat sequence was $(CA)_{16}TA(CA)_{10}$. The AM10 fragment contained two repeat sequences separated by 45 bp of nonrepeat sequence. The first repeat, designated AM10GA, was $(GA)_2ATGACA(GA)_{11}$. The second repeat, designated AM10GT, was not used in this study because upon analysis of the TX-SCA1 kindred it yielded the same information as the AM10GA repeat. The AM10GT repeat consists of $(GA)_2AA(GA)_6GTGA(GT)_{16}AT(GT)_5$. Primer information for AM10GT is available through the Genome Data Base. At SB1, the repeat tract was not sequenced; only flanking sequence was determined.

As there are differences in allele distributions of markers among the different races, allele frequencies are reported here separately for the CEPH kindreds (Caucasian) and the TX-SCA1 kindred (American black) (Table 2). CEPH allele frequencies were based on 72 independent chromosomes for SB1, 82 independent chromosomes for AM10, and on the full set of 40 families for D6S109 and LR40. TX-SCA1 allele frequencies were based on 45 independent chromosomes for LR40, 43 independent chromosomes for SB1, 45 independent chromosomes for AM10, and 42 independent chromosomes for D6S109.

Table 2.
Allele frequencies of new markers

Allele ^a	<u>D6S109^b</u>		<u>AM10GA</u>		<u>SBI</u>		<u>LR40</u>		<u>D6S202^b</u>	
	TXSCA1	CEPH	TXSCA1	CEPH	TXSCA1	CEPH	TXSCA1	CEPH	TXSCA1	CEPH
A ₀	-	0.012	0.070	-	-	-	-	-	-	-
A ₁	0.048	0.024	0.163	0.027	0.244	0.022	0.244	0.022	0.05	0.05
A ₂	0.024	0.220	0.186	0.166	0.045	0.043	0.045	0.043	0.11	0.11
A ₃	0.119	0.024	0.070	0.333	0.111	0.065	0.111	0.065	0.11	0.11
A ₄	0.024	0.232	0.023	-	0.133	0.033	0.133	0.033	0.13	0.13
A ₅	0.071	0.488	0.186	0.097	0.111	0.272	0.111	0.272	0.11	0.11
A ₆	0.261	-	0.093	0.111	-	0.098	-	0.098	0.03	0.03
A ₇	0.024	-	0.093	0.153	0.022	0.054	0.022	0.054	0.22	0.22
A ₈	0.095	-	0.093	0.083	0.045	0.076	0.045	0.076	0.13	0.13
A ₉	0.143	-	-	0.014	0.089	0.054	0.089	0.054	0.08	0.08
A ₁₀	-	-	-	-	0.022	0.065	0.022	0.065	0.03	0.03
A ₁₁	0.048	-	0.023	-	0.133	0.011	0.133	0.011	-	-
A ₁₂	0.048	-	-	-	0.045	0.054	0.045	0.054	-	-
A ₁₃	0.048	-	-	0.014	-	0.097	-	0.097	-	-
A ₁₄	0.071	-	-	-	-	0.033	-	0.033	-	-
A ₁₅	-	-	-	-	-	0.023	-	0.023	-	-

^a Alleles are numbered such that the largest allele is assigned the lowest number and each successive allele is two bp smaller. For D6S109, A₁=215 bp, for Am10GA, A₀=123 bp, for SBI, A₀ = 220 bp, for LR40, TXSCA1 A₁ = 241 bp, (primer set A, Table 1), CEPH A₁ = 267 bp (primer set B, Table 1), for D6S202, A₁ = 154 bp.

^b CEPH data published for D6109 (L.P.W. Ranum et al., *Am. J. Hum. Genet.*, 49, 31-41 (1991) and D6S202 (F. LeBorgne-Demarquoy et al., *Nucl. Acids Res.*, 19, 6060 (1991)).

2. Genetic Linkage Data

a. CEPH families. In order to establish a well-defined genetic map for the SCA1 region, newly isolated DNA markers were mapped using the CEPH reference families. Results of pairwise linkage analyses in CEPH kindreds are shown in Table 3. No recombination was observed between AM10GA and D6S89 ($\theta = 0.00$, $Z_{\max} = 15.1$) using a subset of the CEPH panel which defined 26 recombinants between D6S109 and D6S89. The markers D6S109 and LR40 are close to D6S89, with recombination fractions of 0.067 ($Z_{\max} = 71.4$) and 0.04 ($Z_{\max} = 84.5$) respectively.

Selected multipoint analyses were performed to position the newly isolated markers D6S109, LR40, D6S202 with respect to markers previously mapped using the CEPH panel. The CMAP program was used for three- and four-point linkage analyses to position D6S109 relative to D6S88 and D6S89 and to position LR40 and D6S202 relative to each other and to D6S89 and F13A. For the three-point analyses, the D6S88 - D6S89 interval was fixed based on the two-point recombination fraction in CEPH and the lod score was calculated at various recombination fractions. The order D6S88 - D6S109 - D6S89 is favored over the next most likely order by odds of $4 \times 10^3 : 1$ (Table 4). For the four-point analyses, both the D6S89 - D6S202 - F13A and the D6S89 - LR40 - F13A intervals were fixed based on the two-point recombination fractions; lod scores were then calculated for LR40 and D6S202 at various θ values on the respective fixed maps. The order D6S89 - LR40 - D6S202 - F13A is favored over the next most likely order in both analyses; odds in favor were 400 : 1 when the position of LR40 was varied and were 1×10^6 to 1 when D6S202 was varied (Table 4).

The order of AM10GA and D6S89 could not be determined using the D6S109/D6S89 CEPH recombinants. However, the order AM10GA - D6S89 - SB1 was deduced by characterization of overlapping yeast artificial chromosome clones containing these markers (see below). Furthermore, one end of this contig is present in a well characterized radiation-reduced hybrid known to contain D6S109 and other centromeric markers, indicating the order D6S109 - AM10GA - D6S89 - SB1.

Table 3.
Pairwise linkage results in CEPH

Marker Pair	$\theta_m = \theta_r$	Z_{\max}	θ_m	θ_r	Z_{\max}	χ^2
HLA and D6S88	0.128	26.4	0.103	0.168	26.8	1.86
D6S109	0.126	48.4	0.062	0.176	51.0	12.1*
AM10	0.608	0.0440	0.301	0.500	0.246	0.929
D6S89	0.158	43.3	0.091	0.225	46.6	15.2*
SB1	0.574	0.0190	0.299	0.500	0.400	0.381
LR40	0.213	25.5	0.116	0.306	30.0	20.8*
HZ30	0.251	21.6	0.191	0.318	23.6	8.95*
F13A	0.291	8.81	0.255	0.326	9.14	1.52
D6S88 and D6S109	0.017	48.6	0.024	0.009	48.8	0.846
AM10	0.654	0.0290	0.499	0.696	0.047	0.0820
D6S89	0.086	36.1	0.076	0.098	36.2	0.0750
SB1	0.203	1.09	0.136	0.687	1.36	1.27
LR40	0.088	31.1	0.078	0.104	31.2	0.350
HZ30	0.135	30.4	0.124	0.152	30.4	0.340
F13A	0.180	10.2	0.158	0.217	10.3	0.626
D6S109 and AM10	0.730	0.933	0.170	0.502	1.67	3.39
D6S89	0.067	71.4	0.035	0.090	72.5	5.15*
SB1	0.742	1.95	0.113	0.501	4.32	10.9*
LR40	0.109	50.6	0.050	0.152	52.9	10.5*
HZ30	0.162	36.6	0.147	0.174	36.7	0.515
F13A	0.207	14.4	0.211	0.204	14.4	0.0368
AM10 and D6S89	0.000	15.1	0.000	0.000	15.1	0.000
SB1	0.000	13.2	0.000	0.000	13.2	0.000
LR40	0.021	8.74	0.000	0.050	9.11	1.74
HZ30	0.000	13.8	0.000	0.000	13.8	0.000
F13A	0.135	3.48	0.042	0.253	4.39	4.16*

-45-

D6S89 and SB1	0.000	25.0	0.000	0.000	25.0	0.000
LR40	0.040	84.5	0.030	0.049	84.7	0.925
HZ30	0.078	76.0	0.075	0.077	76.0	0.0230
F13A	0.151	30.7	0.139	0.160	30.7	0.248
SB1 and LR40	0.033	14.4	0.022	0.044	14.5	0.350
HZ30	0.026	17.5	0.032	0.020	17.5	0.0300
F13A	0.136	4.80	0.119	0.155	4.84	0.170
LR40 and HZ30	0.079	64.8	0.092	0.050	65.0	1.09
F13A	0.131	29.1	0.121	0.140	29.2	0.189
HZ30 and F13A	0.109	38.4	0.122	0.106	38.4	0.0092

*Indicates statistically significant differences were observed in the recombination fractions when the assumption of homogeneity ($\theta_m = \theta_f$) was rejected; that is the likelihood that $\chi^2 > 3.84$ with 1 degree of freedom should occur by chance in $P < 0.05$.

-46-

Table 4.
Three and four point linkage analyses in the CEPH families

<u>Order</u>	<u>Z_{max}</u>	<u>Relative Odds</u>	<u>Odds in favor</u>
D6S109-D6S88-D6S89	90.6	2X10 ⁸	
D6S88-D6S109-D6S89	94.2	8X10 ¹¹	4X10 ³
D6S88-D6S89-D6S109	82.3	1	
LR40-D6S89-D6S202-F13A	96.1	1X10 ³⁴	
D6S89-LR40-D6S202-F13A	98.6	4X10 ³⁶	400:1
D6S89-D6S202-LR40-F13A	73.9	8X10 ¹¹	
D6S89-D6S202-F13A-LR40	62.0	1	
D6S202-D6S89-LR40-F13A	89.5	1X10 ³²	
D6S89-D6S202-LR40-F13A	57.5	1	
D6S89-LR40-D6S202-F13A	95.5	1X10 ³⁸	10 ⁶ :1
D6S89-LR40-F13A-D6S202	77.6	1X10 ²⁰	

-47-

b. SCA1 kindreds. Results of pairwise linkage analyses in SCA1 kindreds are shown in Table 5. AM10GA, D6S89, and SB1 are all closely linked to SCA1. No recombination was observed between AM10GA and SCA1; the lod score is 42.1 at a recombination fraction of 0.00. The recombination fraction
5 between D6S89 and SCA1 is 0.004 (lod score of 67.6). The recombination fraction between SB1 and SCA1 is 0.007 (lod score of 39.5). D6S109, LR40 and D6S202 are linked to SCA1 as well, but at greater distances (recombination fractions of 0.04, 0.03, and 0.08 respectively). Based on genetic mapping in nine large kindreds, the SCA1 locus is very close to D6S89 and AM10GA, with a $Z_{\max}-1$ support interval
10 less than or equal to 0.02 in both cases.

Table 5.
Pairwise lod scores for SCA1 and dinucleotide repeat markers

	0	Recombination fraction							Support θ^a	Interval ^b
		0.001	0.05	0.1	0.2	0.3	0.4	Z^a		
SCA1:D6S109	-∞	22.68	33.81	32.03	25.19	16.56	7.24	33.82	0.04	0.02 to 0.09
SCA1:AM10GA	42.14	42.06	38.48	34.51	25.86	16.63	7.30	42.14	0.00	0.00 to 0.02
SCA1:D6S89	-∞	67.35	62.78	56.39	42.51	27.56	12.09	67.58	0.004	0.00 to 0.02
SCA1:SB	-∞	39.02	37.33	33.92	26.16	17.53	8.33	39.46	0.007	0.00 to 0.03
SCA1:LR40	-∞	27.80	31.77	29.73	23.61	16.11	7.77	32.08	0.03	0.001 to 0.07
SCA1:D6S202	-∞	4.41	25.80	26.47	22.12	14.77	6.51	26.61	0.08	0.04 to 0.14

^a Z = maximum lod score, θ = recombination fraction at maximum lod score.

^b Z_{\max}^{-1} = support interval for θ (Cytogenet Cell Genet, 40, 356-359 (1985)).

3. Analysis of Key Recombinants

One recombination event between D6S89 and SCA1 has been confirmed in an affected individual. The patient, individual MI-2 in Figure 4, was also recombinant at SB1, although uninformative at LR40 and D6S202. He carried a disease haplotype at the HLA, D6S109 and AM10 loci, demonstrating that SCA1 is centromeric to D6S89, as indicated by the rightmost arrow in Figure 4. To eliminate the possibility of sample mix-up, the patient's DNA was reextracted from a hair sample and retyped for D6S109, D6S89, D6S202, LR40, AM10GA, and SB1. The results from the hair sample matched those from the cell line originally established from the patient's blood. The patient's medical records were carefully reexamined and it was confirmed that he did indeed have ataxia. In addition, his haplotypes were consistent with those of a sister and a daughter.

D6S109 lies centromeric to D6S89; six recombination events have been observed between D6S109 and SCA1, as shown in Figure 4. At this point, D6S109 is the centromeric marker closest to SCA1. The arrows in Figure 4 denote the maximum region common to all affected chromosomes, and therefore the maximum possible region containing the SCA1 gene, which extends from D6S89 to D6S109.

No additional marker-SCA1 recombination events have been observed between D6S89 and SB1. Markers further telomeric to SB1 show additional recombination with SCA1 -- one recombination event between SCA1 and LR40 and three recombination events between SCA1 and D6S202. These events are depicted in Figure 4 (all recombination events depicted in Figure 4 are in affected individuals).

II. Mapping and Cloning the Critical Region for the SCA1 Gene

A 2.5-Mb yeast artificial chromosome (YAC) contig was developed with the ultimate goal of defining and cloning the region likely to contain the SCA1 gene (SCA1 critical region).

A. Materials and Methods

1. Cell lines

I-7 is a human-hamster hybrid cell line which contains the short arm of chromosome 6 as its only human chromosome. See, H.Y. Zoghbi et al., Genomics, 6, 352-357 (1990). R86, R78, R72, R54 and R17 are radiation reduced hybrid cell lines retaining various portions of 6p22-p23. See, H.Y. Zoghbi et al., Genomics, 9, 713-720 (1991). R54 retains markers known to be telomeric to D6S89, such as D6S202 and F13A.

2. Generation of new DNA markers and Sequence-Tagged Sites (STSs)

DNA from a radiation reduced hybrid retaining D6S89 (R86) and DNAs from four radiation hybrids (R78, R72, R54 and R17) which do not retain D6S89 but retain markers immediately flanking D6S89 were used in comparative *Alu*-PCR to isolate region-specific DNA markers. See, D.L. Nelson et al., Proc. Natl. Acad. Sci. USA, 86, 6686-6690 (1989); and H.Y. Zoghbi et al., Genomics, 9, 713-720 (1991). In addition, R78 was useful in eliminating markers derived from the centromeric region of 6p. H.Y. Zoghbi et al., Genomics, 9, 713-720 (1991). *Alu*-PCR was carried out using *Alu* primers 559 and 517 individually (D.L. Nelson et al., Proc. Natl. Acad. Sci. USA, 86, 6686-6690 (1989)) as well as PDJ 34 (C. Breukel et al., Nucleic Acids Res., 18, 3097 (1990)). *Alu*-PCR fragments found to be present in R86 but absent in R78, R72, R54 and R17 were identified and were cloned into *EcoRV*-digested pBluescript IKS+ plasmid (Stratagene, La Jolla, CA) which was modified using the T-vector protocol. See, D. Marchuk et al., Nucleic Acids Res., 19, 1154 (1990). Cloned fragments were sequenced on an Applied Biosystems, Inc. (Foster City, CA) automated sequencer to establish STSs.

3. Isolation and Characterization of YAC clones

The Washington University YAC library (B.H. Brownstein et al., Science, 244, 1348-1351 (1989)), and the CEPH YAC library (H.M. Albertsen, et al., Proc. Natl. Acad. Sci. USA, 87, 4256-4260 (1990)), were screened using a PCR-based method. See, E.D. Green et al., Proc. Natl. Acad. Sci. USA, 87, 1213-1217 (1990); and T.J. Kwiatkowski et al., Nucleic Acids Res., 18, 7191-7192 (1990). PCR amplifications were carried out in 25-50 ml final volume with 50 mM KCl, 10

-51-

mM Tris-HCl pH 8.3, 1.25 mM MgCl₂, 0.01% (w/v) gelatin, 250 μM of each dNTP; 1.25 units of Amplitaq polymerase (Perkin-Elmer, Norwalk, CT) and 1 μM of each primer. PCR cycle conditions are specified in Table 6.

Table 6.
STSs and YACs in 6p22-p23

Probe	Primer set	YACs ^a	Annealing temp. ^b
D6S89	cttgttcacatgccttgtgcaccta agcgactgcctaaac	B126G2, B134D5, B172B3, B214D3, C5C12, 191D8, 299B3, 379C2, 468D12, 124G2, 511H11	55°C
AM10 (D6S335)	ttaaggaagtgttcacatcaggg aattgtgcttatgtcactggg	A23C3, A183C6, A250D5, B238F12, A91D2	55°C
A250D5-L (D6S337)	aattctggagagaggatgttggt tcctttttgtag	195B5, 242C5, 475A6, 30F12	44°C
64U	catcgtgttgtgttggaagctc agacgctaaactcaagg	492H3, 172B5, 227B1, 261H7	50°C
D6S288	atgatccgtggtagtggcagga cctgttactgacgcc	60H7, 351B10	55°C
D6S274	ctcatctgttgaatgggatctta aatgctatgccttccg	486F9, 149H3, 42A5, 283B2, 320E12	55°C
FLB1 (D6S339)	tgcaaatccctcagttcacttgctt gactttgccatgttc	140H2, 270D3, 274D12, 401D6, 57G3, 168F1	50°C
AM12 (D6S336)	ataccatacggattgagggca acactatcaggctaagaatg	A71B3, 228A1, 193B3, 90A12, 539C11, 53G12, 35E8	55°C
53G12-L	caaataccagcaactcaccagc gggtccttcagcatcctacattc	3G6, 82G12, 98G5, 135F6, 198C8, 330G1	58°C

^a YACs in this study are from the CEPH and Washington University libraries. I.D. numbers identify the library source (Washington University I.D. numbers are preceded by a letter). Several YACs were identified with more than one STS; for such information, please refer to Table 2.

^b PCR conditions were 94°C for 4 minutes followed by 35-40 cycles of 94°C denaturation for 1 minute, annealing at the specified temperature for 1 minute, and 72°C extension for 2 minutes. A final extension step of 7 minutes at 72°C was used. PCR buffer and primer concentrations are as described in the text; for the 53G12-L STS a final concentration of 2% formamide was used in the PCR reaction.

Yeast DNA-agarose blocks were prepared as described by D.C. Schwartz et al., Cell, **37**, 67-75 (1984); and G.J.B. van Ommen et al. in Human Genetic Diseases-A Practical Approach; K.E. Davies, ed.; pp. 113-117; IRL Press, Oxford (1986). All the YAC clones were analyzed by pulsed-field gel electrophoresis (PFGE) to determine the insert size and to confirm that a single YAC was present in a specific colony. YAC inserts were sized by electrophoresing yeast DNA through a 1% Fastlane agarose (FMC, Rockland, ME) gel in 0.5x TAE (20 mM Tris-acetate/0.5 mM EDTA). For rapid detection of possible overlaps between YAC clones isolated at different STSs, the labelled *Alu*-PCR products of new YACs were hybridized to filters containing *Alu*-PCR products of individual YACs in the region. Most of the YAC clones were tested for chimerism using the *Alu*-PCR dot blot method described by S. Banfi et al., Nucleic Acids Res., **20**, 1814 (1992). The *Alu*-PCR products from YAC clones were hybridized to a dot-blot containing the *Alu*-PCR products from monochromosomal or highly reduced hybrids representing each of the 24 different human chromosomes as previously described by S. Banfi et al., Nucleic Acids Res., **20**, 1814 (1992). In addition a dot-blot containing *Alu*-PCR products from radiation reduced hybrids representing different segments of 6p was used to insure that a YAC does not contain two non-contiguous segments from 6p. Ends of YAC clones were isolated either by inverse-PCR as previously described by G. Joslyn et al., Cell, **66**, 601-613 (1991) or by *Alu*-vector PCR as described by D.L. Nelson et al., Proc. Natl. Acad. Sci. USA, **88**, 6157-6161 (1991). *Alu*-vector PCR was carried out using *Alu*-primers PDJ34 and SAL1, as described by C. Breukel et al., Nucleic Acids Res., **18**, 3097 (1990); and the pYAC4 vector primers described by M.C. Wapenaar et al., Hum. Mol. Genet., **2**, 947-995 (1993) and analogous vectors described by G.P. Bates et al., Nature Genetics, **1**, 180-187 (1992). All YAC ends were regionally mapped by hybridization to Southern blots containing *Eco*RI-digested DNAs from the YAC clones and from the hybrid cell lines: I-7, R86, and R72.

30 4. Cosmid library preparation from YACs

Cosmid libraries were prepared from four YAC clones; 227B1, 195B5, A250D5, and 379C2. Genomic DNA from YACs was partially digested with *Mbo*I and cloned into cosmid vector superCos 1 (Stratagene, La Jolla, CA) following the

manufacturer's recommendations. Clones containing human inserts were identified using radiolabeled sheared human DNA as a probe.

5. Long range restriction analysis

5 YAC plugs were digested to completion using rare-cutter restriction enzymes as described by M.C. Wapenaar et al., Hum. Mol. Genet., **2**, 947-995 (1993) and analogously by G.A. Silverman et al., Proc. Natl. Acad. Sci. USA, **86**, 7485-7489 (1989). Enzymes were purchased from New England Biolabs (Beverly, MA) and Boehringer Mannheim Biochemicals (Indianapolis, IN) and were used as
10 recommended by the manufacturer. All PFGE analyses were performed on a Bio-Rad CHEF apparatus under conditions that separate DNA fragments in the 50 kb to 600 kb range. The gels were stained with ethidium bromide, and either acid nicked or subjected to 200,000 mJ of UV energy in a UV Stratalinker 1800 (Stratagene, La Jolla, CA). The gels were denatured in 0.4 N NaOH and transferred to Sure Blot
15 hybridization membrane (Oncor, Gaithersburg, MD) in either 10xSSC (1.5 M NaCl/150 mM NaCitrate) or 0.4 N NaOH according to the manufacturer's recommendations. Hybridizations of the filters were carried out using the probes listed in Table 6 and Figure 6. Also pBR322 *Bam*HI/*Pvu*II fragments of 2.5 kb and 1.6 kb specific for the left (TRP/CEN) and right (URA) pYAC4 vector arms
20 respectively, were used. Probes were radiolabelled using the random priming technique described by A.P. Feinberg et al., Anal. Biochem., **137**, 266-267 (1984); repetitive sequences were blocked using sheared human placental DNA as previously described by P.G. Sealy et al., Nucleic Acids Res., **13**, 1905-1922 (1985).

25 6. Dinucleotide repeat analysis

Primer sequences and PCR cycle conditions are presented in Table 6. Buffer conditions were the same as for *Alu*-PCR. All reaction volumes were 25 µl and contained 40 ng of genomic DNA. One primer of each pair was labelled at the 5' end with [γ -³²P] dATP. Four microliters of each reaction was mixed with 2 µl
30 formamide loading buffer, denatured at 90-100°C for 3 minutes, cooled on ice and 4-6 µl was used for electrophoresis on a 4% polyacrylamide/7.65 M urea sequencing gel.

B. Results

1. Generation of sequence tagged sites in 6p22-p23 and YAC screening

Comparative analysis of the *Alu*-PCR products from the radiation hybrid, which retains D6S89 (R86) and from the four radiation hybrids deleted for D6S89 but retaining markers which flank D6S89 (R78, R72, R54 and R17) allowed the identification of three new DNA fragments that were present in R86 but absent in the other four. These three DNA fragments termed, AM10, AM12 and FLB1 were isolated and mapped using a 6p somatic cell hybrid panel and the radiation reduced hybrid panel (H.Y. Zoghbi et al., Genomics, 9, 713-720 (1991)) to confirm their regional localization. All three mapped to 6p and to R86 confirming their close proximity to the D6S89 locus. These three *Alu*-PCR fragments were subcloned and sequenced to establish sequenced tagged sites (STSs). STSs at AM10, AM12, FLB1 and D6S89 were used to screen the Washington University and the CEPH YAC libraries (H.M. Albertsen, et al., Proc. Natl. Acad. Sci. USA, 87, 4256-4260 (1990); and B.H. Brownstein et al., Science, 244, 1348-1351 (1989)). YACs isolated at these four STSs were analyzed for overlap. Insert termini from the YACs representing contig ends were isolated, subcloned and were sequenced to establish new STSs for further YAC walking. In one case an STS was established by using a subclone from a cosmid derived from a cosmid library generated for YAC 195B5.

Recently several highly informative dinucleotide repeat markers have been identified and mapped genetically by J. Weissenbach et al., Nature, 359 794-801 (1992). As discussed above, two markers, D6S274 and D6S288 were found to map within the SCA1 critical region and were subsequently used to screen the YAC libraries. Using the STSs listed in Table 6, YAC clones were isolated.

2. Characterization of YAC clones

The sizes of the YAC inserts were determined by pulsed-field gel electrophoresis (PFGE); insert sizes ranged from 75-850 kb. Given the high frequency of insert chimerism, an *Alu*-PCR based hybridization strategy for rapid detection of chimerism, as described by S. Banfi et al., Nucleic Acids Res., 20, 1814 (1992) was used. Thirty of the YAC clones were tested using this approach and eight (27%) were found to be chimeric. Insert ends isolated from YACs determined

to be non-chimeric by the dot blot hybridization approach mapped to 6p22-p23 with the exception of the two ends from 198C8 which proved to map to other chromosomes.

Two approaches were used, inverse-PCR (G. Joslyn et al., Cell, 66, 601-613 (1991)) and *Alu*-PCR (analogous to that described by D.L. Nelson et al., Proc. Natl. Acad. Sci. USA, 86, 6686-6690 (1989)) to isolate YAC ends. In total, 34 YAC ends were isolated; inverse-PCR yielded 26 ends and *Alu*-vector PCR yielded 8 ends. To isolate the left end of the 195B5 YAC we screened a cosmid library prepared from this YAC using pYAC4 left end sequences (S.K. Bronson et al., Proc. Natl. Acad. Sci. USA, 88, 1676-1680 (1991)) as a probe. This approach was taken because inverse-PCR yielded an end which was predominantly an *Alu*-containing sequence and *Alu*-PCR failed in yielding an end. Cosmid clone A32 was found to contain the left end of 195B5 and a subclone, 64U, was used to establish an STS for further YAC library screenings.

In order to confirm the 6p22-p23 regional origin of all YAC ends or subclones, these fragments were used as probes against Southern blots containing *Eco*RI-digested DNAs from a somatic cell hybrid retaining 6p (I-7), from radiation reduced hybrids known to retain fragments of 6p (H.Y. Zoghbi et al., Genomics, 9, 713-720 (1991)) and from the YAC clones at a particular STS.

3. Probe content mapping of YACs

In order to define the degree of overlap between the clones and to detect possible rearrangements such as internal deletions of the YACs, a probe content mapping strategy was used based on: 1) PCR analysis of all the clones using all the STSs in the region including both the ones described in Table 6, and those at highly informative dinucleotide repeats such as AM10-GA and SB1; and 2) hybridization of Southern blots containing *Eco*RI-digested DNAs from YACs in the relevant region, with densely-spaced DNA probes derived from YAC ends, cosmids subclones of YACs, or *Alu*-PCR fragments from YACs. The results of this analysis for a representative subset of the YACs (32 clones) are summarized in Table 7. Thirty-nine YAC clones form an uninterrupted YAC contig from D6S274 to 82G12-R (right end of YAC clone 82G12). Other than an internal deletion in one YAC (351B10) no other deletions were detected within the resolution of this analysis;

-57-

furthermore the extent of chimerism for some YAC clones (such as 270D12 and 140H2) was determined. The centromere-telomere orientation of the YAC contig on 6p was determined using both genetic data as well as physical mapping data. Using dinucleotide repeats analysis at D6S109, AM10GA, D6S89, and SB1 in the
5 key individual with recombination event between D6S89 and SCA1 revealed that the recombination event occurred between AM10GA and D6S89. Given that D6S109 is centromeric to D6S89, the recombination analysis suggests that AM10GA is centromeric to D6S89. The centromere-telomere position of SB1 with respect to D6S89 could not be determined genetically.

TABLE 7.**Characterization of YACs using 6p22-p23 STSs and YAC fragments**

YAC	Size (kb)	Chimerism	D6S274	60H7Lg	D6S288	64U	A25005-L	AM10-GA	AM10	168F1-R	C5C12-R	D6S89	B214D3-R	FLB1	53G12-R	401D6-R	AM12	135F6-L	53G12-L	135F6-R	83G12-R
149H3	345	N	+	+	-	-	-														
60H7	580	N	+	+	+	-	-														
351B10	330	N	+	-	+	-	-														
227B1	560	N	+	+	+	+	-														
172B5	345	Y	-	-	+	+	-														
195B5	365	N	-		-	+	+														
475A6	365	N				-	+														
242C5	340	N				-	+	+													
A250D5	250	N				-	+	+													
A23C3	530	Y				-	-	+													
A18306	120	N				-	-	-	+												
B238F12	390	Y				-	-	+	+												
A91D2	325	N				-	-	-	+												
191D8	650	N				-	-	-	+	+											
379C2	575	N				-	-	-	+	+											
C5C12	75	N				-	-	-	+	+											
B214D3	200	N				-	-	-	+	+											
299B3	375	N				-	-	-	+	+											
468D12	280	N				-	-	-	+	+											
168F1	400	N				-	-	-	+	+											
270D3	650	Y				-	-	-	+	+					+	+					
274D12	240	N				-	-	-	+	+					+	+					
140H2	440	Y				-	-	-	+	+					+	+					
57G3	400	N				-	-	-	+	+					+	+					
401D6	340	N				-	-	-	+	+					+	+					
193B3	850	Y				-	-	-	+	+					+	+					
228A1	350	Y				-	-	-	+	+					+	+					
90A12	650	Y				-	-	-	+	+					+	+					
35E8	400	N				-	-	-	+	+					+	+					
53G12	370	N				-	-	-	+	+					+	+					
135F6	400	N				-	-	-	+	+					+	+					
82G12	380	N				-	-	-	+	+					+	+					

Note. (+) = present, (-) = absent; Y/N = chimerism is/not detected. YAC ends are identified by YAC names followed by L or R for left or right.

Physical mapping, using both radiation hybrids and YACs, was carried out to resolve the centromere-telomere order of the loci. The radiation reduced hybrids R17 and R72 are known to contain markers centromeric to D6S89; these markers include D6S108 and D6S88 which map centromeric to D6S109. See, H.Y. Zoghbi et al., Genomics, 2, 713-720 (1991). R72 also retains D6S109, but a small gap in R17 was revealed as this radiation hybrid did not retain D6S109, but was positive for an end isolated from a YAC at the D6S109 locus. Analysis of the radiation reduced hybrids revealed that D6S274 and D6S288 are present in R17, R72 and R86, whereas AM10GA, D6S89, and SB1 are present only in R86 (Figure 5). Furthermore, STS content mapping with D6S260 and D6S289, two dinucleotide repeats that are telomeric to D6S288 (J. Weissenbach et al., Nature, 359 794-801 (1992)), revealed that D6S260 is present in the same YACs as D6S89 and SB1 (379C2 and 168F1), and that D6S289 is present in 57G3 and 35E8 two YACs derived using the FLB1 and AM12 STS respectively. These data, confirm that the order of the loci as well as the centromere-telomere orientation of the YAC contig presented in Figure 6 is correct.

Figure 6 shows a selected subset of YAC clones which span the entire contig from D6S274 to 82G12-R. A minimal number of 8 YACs spans this region. The positions of the STSs which were used to isolate the YACs are also shown. Based on the size of the YACs and the degree of overlap, this contig is estimated to span 2.5 Mb of genomic DNA in 6p22-p23 with D6S89 located approximately in the middle.

4. Delineating the SCA1 critical region

Genetic studies using recently identified dinucleotide repeats (AM10GA and SB1) showed that SCA1 maps centromeric to the D6S89 locus very close to AM10GA (peak load score of 42.1 at a recombination frequency of zero) in nine large SCA1 kindreds (Example 1, above). Thus D6S89 is the closest flanking marker at the telomeric end. Previously, the closest flanking marker at the centromeric end was D6S109, a dinucleotide repeat estimated to be 6.7 cM centromeric to D6S89. To identify a closer flanking marker at the centromeric end, we mapped D6S260, D6S274, D6S288 and D6S289, four dinucleotide repeat-containing markers known to map 6p22-p23 (J. Weissenbach et al., Nature, 359

794-801 (1992)). The regional mapping of these markers was done using radiation reduced hybrids and the YAC clones isolated from this region. These data revealed that D6S274 and D6S288 map centromeric to AM10GA as evident by amplification of DNA from radiation hybrids R17 and R72 which are known to be centromeric to AM10GA. Genotypical analysis of the DNAs from individuals with key recombination events between D6S109 and D6S89 as well as from affected and normal individuals (to establish chromosomal phase) from the five SCA1 kindreds (MN-SCA1, MI-SCA1, TX-SCA1, M-SCA1 and MS-SCA) was carried out. This analysis revealed no recombination between D6S288 and SCA1. A single recombination event between D6S274 and D6S288 was detected in individual MN-1 from the MN-SCA1 kindred (Figure 7); this individual was one of the six individuals identified above as having a recombination event between SCA1 and D6S109. This analysis allowed us to identify D6S274 as the closest flanking marker at the centromeric end. These data combined with that discussed above determined that the SCA1 critical region maps between D6S274 and D6S89. This candidate region (1.2 Mb) is cloned in a minimum of four overlapping and non-chimeric YACs as shown in Figure 8.

5. Long-range restriction mapping

In order to have an estimate of the size of the YAC contig in the SCA1 critical region we performed long-range restriction analysis on YACs from this region. The YACs used for this analysis included: 227B1, 60H7, 351B10, 172B5, 195B5, A250D5, 379C2, and 168F1. The following rare-cutter restriction enzymes were used: *NotI*, *BssHII*, *NruI*, *MluI*, and *SacII*. Restriction fragments separated by PFGE and transferred onto nylon membranes, were detected by sequential hybridizations of the filter to several DNA probes which included: DNA probes specific for the left and right arm of the pYAC4 vector; insert termini for internal YAC clones; internal probes and cosmid subclones; and an *Alu*-specific probe. The position and names of all the probes used in the long-range restriction analysis is shown in Figure 8. Based on this analysis the internal deletion for YAC 351B10 was confirmed. The extent of overlap between the YAC clones was determined. The size of the critical SCA1 region was estimated to be 1.2 Mb. Internal deletions and/or other rearrangements could not be excluded for the areas where a single YAC

was analyzed by restriction enzyme analysis. These include approximately a 220 kb region within YAC 195B5 and a 335 kb region within YAC 379C2.

III. Expansion of an Unstable Trinucleotide Repeat in SCA1

A. Methods

1. Screening for trinucleotide repeats

Genomic DNA from YACs was partially digested with *Mbo*I and cloned into cosmid vector super CosI (Stratagene) following the manufacturer's protocol. Clones containing human inserts were identified by hybridization with radiolabeled human DNA and were arrayed on a gridded plate. Filter lifts of cosmid clones from YAC227B1 were screened for the presence of trinucleotide repeats by hybridization to [γ -³²P] end-labelled (GCT)_n oligonucleotide. In a parallel experiment, a mixture of 10 oligonucleotides representing the various permutations of trinucleotide repeats were end-labelled and hybridized to a Southern transfer of *Eco*RI-digested cosmids from YACs 195B5 and A250D5. Hybridizations were done in a solution of 1 M NaCl, 1% sodium dodecyl sulfate (SDS) and 10% (w/v) dextran sulphate. Filters were washed in 2xSSC (1xSSC is 0.15 M sodium chloride and 0.015 M sodium citrate), and 0.1% SDS at room temperature for 15 minutes, followed by a 15 minute wash at room temperature in a solution prewarmed to 67°C. Both strategies identified several positive clones, 22 of which were overlapping and contained the same 3.36-kb *Eco*RI fragment which hybridized to the (GCT)_n probe and ultimately proved to have the CAG repeat by sequence analysis.

2. Genomic digests and Southern blots

Genomic DNAs were digested with *Taq*I (Boehringer Mannheim, Indianapolis, IN) or *Bst*NI (New England Biolabs, Beverly, MA) according to the manufacturers recommendations. Southern blotting was done following standard protocols.

3. DNA sequencing

To determine the DNA sequence in the region containing and flanking the CAG trinucleotide repeats, clone pGCT-7, containing the 3.36 kb-*Eco*RI

-62-

fragment, was subcloned. A 400-bp fragment with CAG trinucleotide repeats was generated from pGCT-7 by *Sau3AI* digestion and subcloned into the *Bam*HI site of pBluescriptKS- (Stratagene, La Jolla, CA) (clone pGCT-7.s1). In addition, pGCT-7 was digested with *Pst*I to remove 1.3 kb of DNA and recircularized for transformation (clone pGCT-7.p2). The position of the trinucleotide repeats was determined by PCR using (GCT)₇ oligonucleotide and one of the flanking sequencing primers as PCR primers. Initial results indicated that the CAG trinucleotide repeats were on the reverse primer strand, about 1.3 kb from the reverse primer, that is, 400 bp from the *Pst*I site. DNA sequencing was performed by di-deoxynucleotide chain-termination method using Sequenase and Δ Taq Cycle-Sequencing kit (United States Biochemical, Cleveland, OH). Both universal (-40) and reverse primers were used for clone pGCT-7.s1, while only universal (-40) primer was used for sequencing pGCT-7.p2.

4. RT-PCR and Northern analysis

Total RNA was extracted from lymphoblastoid cells using guanidinium thiocyanate followed by centrifugation in a cesium chloride gradient. Poly(A)⁺ RNA was selected using Dynabeads oligo(dT)₂₅ from Dynal (Great Neck, NY). First strand cDNA synthesis was carried out using MMLV reverse transcriptase (BRL, Gaithersburg, MD). RT-PCR was carried out using hot start PCR with three cycles of: 97°C for 1 minute, 59°C for 1 minute, and 72°C for 1 minute for the Pre1 and Pre2 primer set. Following that 33 cycles of 94°C for 1 minute, 57°C for 1 minute, and 72°C for 1 minute were carried out. For the Rep1 and Rep2 primer pair the same PCR cycling conditions were followed at lower annealing temperatures of 57°C and 55°C respectively. The RT-PCR products were analyzed on 6% Nusieve agarose gel. The northern blot containing various human tissues was purchased from Clontech (Palo Alto, CA).

5. PCR Analysis

Fifty ng of genomic DNA was mixed with 5 pmol of each primer (CAG-a/GAG-b or Rep-1/Rep-2) in a total volume of 20 μ l containing 1.5 mM MgCl₂, 300 μ M dNTPs (1.25 mM MgCl₂ and 250 μ M dNTPs for Rep-1/Rep-2 primers), 50 mM KCl, 10mM Tris-HCl pH 8.3, and 1 unit of Amplitaq (Perkin

-63-

Elmer, Norwalk, CT). For the CAG-a/CAG-b primer pair [α - 32 P]dCTP was incorporated in the PCR reaction, for Rep-1/Rep-2 primer pair the Rep-1 primer was labeled at the 5' end with [γ - 32 P]dATP. Formamide was used at a final concentration of 2% when using the Rep-1/Rep-2 primer pair. Samples, overlaid with mineral oil, were denatured at 94°C for 4 minutes followed by 30 cycles of denaturation (94°C, 1 minute), annealing (55°C, 1 minute), and extension (72°C, 2 minutes). Six microliters (μ l) of each PCR reaction was mixed with 4 μ l formamide loading buffer, denatured at 90°C for 2 minutes, and electrophoresed through a 6% polyacrylamide/7.65 M urea DNA sequencing gel. Allele sizes were determined by comparing migration relative to an M13 sequencing ladder.

B. Results

1. Cloning of the CAG repeat region in SCA1

As discussed above, in efforts to clone the SCA1 gene, key recombination events were analyzed using several dinucleotide repeat polymorphisms mapping to 6p22-p23 to identify the minimal region likely to contain the SCA1 gene. This analysis revealed that there were no recombination events between SCA1 and the centromeric marker D6S288 in five large kindreds or between SCA1 and the telomeric marker AM10GA in nine large kindreds. A single recombination event was detected between D6S274 and D6S288 identifying the closest flanking marker at the centromeric end to be D6S274. At the telomeric end, a single recombination event was detected between AM10GA and D6S89 and identified the latter as the flanking marker. A yeast artificial chromosome (YAC) contig extending from D6S274 to D6S89 and spanning the entire SCA1 candidate region was developed. A subset of the YAC clones encompassing this region is shown in Figure 9. Long-range restriction analysis determined the size of the SCA1 candidate region to be approximately 1.2 Mb. Cosmid libraries were constructed from YACs 227B1, 195B5, A250D5, and 379C2. Arrays of cosmid clones containing human inserts were hybridized with an oligonucleotide consisting of tandemly repeated CAG, as well as with oligonucleotides containing other trinucleotide repeats. Several hybridizing cosmid clones were identified, 23 of which were positive for the CAG repeat and mapped to the region between D6S288

and AM10GA (Figure 9). All 22 of these clones shared a common 3.36-kb *EcoRI* fragment that specifically hybridized to the CAG repeat.

2. Variability of the CAG Repeat Using Southern Analysis

5 To test the genetic stability of this repeat in SCA1, we used Southern blotting analysis to examine families with juvenile onset SCA1. A two-generation reduced pedigree from the TX-SCA1 family is shown in Figure 10a. Paternal transmission of SCA1 with an expansion of a *TaqI* fragment was noted. A 2830-bp fragment was detected in DNA from the unaffected spouse and on the normal
10 chromosome from SCA1 patients, whereas a 2930-bp fragment was found in DNA from the affected father (onset at 25 years) and a 3000-bp fragment was detected in DNA from his affected child with an onset at 4 years. In a second SCA1 kindred, family MN-SCA1 (Figure 10b), two offspring inherited SCA1 from their father and differed in their age at onset (25 years and 9 years). These individuals also differ in
15 the size of the amplified *TaqI* fragment they inherited from their affected father, 2900-bp and 2970-bp, respectively.

 Enlargement of the (CAG)_n-containing fragment on SCA1 chromosomes from the same TX-SCA1 juvenile onset family was also demonstrated by Southern analysis following *BstNI* digestion. The *BstNI* fragment is 530-bp on
20 normal chromosomes, is 610-bp in the SCA1 affected father, and is 680-bp in the affected juvenile onset offspring (Figure 10c). In each of these families, nonpaternity was excluded by genotypic analysis with a large number (greater than 10) of dinucleotide repeat markers. In addition, the size of the (CAG)_n-containing *TaqI* fragment in DNA from 30 unaffected spouses was compared to the sizes of the
25 repeat containing *TaqI* fragment in DNA from 62 individuals affected with late-onset SCA1. The affected individuals are from five different SCA1 families: LA-SCA1, MI-SCA1, MN-SCA1, MS-SCA1, and TX-SCA1. In all 30 unaffected spouses fragment sizes were approximately 2830-bp and no expansions or reductions were detected with transmission to offspring. In contrast, DNA from 58
30 of the 62 SCA1 affected individuals contained detectably expanded *TaqI* fragments ranging in size from 2860-bp to 3000-bp in addition to the 2830-bp fragment. The DNAs from the remaining four individuals were found to have an expansion when analyzed by polymerase chain reaction (PCR). The expanded fragment always

segregated with disease, and in some cases the fragment expanded further in successive generations. In the juvenile cases the expanded restriction fragment was larger than that in the affected parent (uniformly the father in the cases analyzed) supporting the conclusion that a DNA sequence expansion is the mutational basis of SCA1.

3. Genomic DNA analysis of repeat regions

To identify the region involved in the DNA expansion, a 500-bp (CAG)_n-containing subclone of the 3.36-kb *Eco*RI fragment was sequenced, as was the entire 3.36-kb fragment (Figure 1). This normal allele demonstrated 30 CAG repeat units. In two of the repeat units (position 13 and 15) a T was present instead of a G.

The expansion of the trinucleotide repeat was observed in all affected individuals examined by PCR from five different kindreds representing at least two ethnic backgrounds, American Black and Caucasian. Genotypic analysis using DNA markers that are very closely linked to SCA1 (D6S274, D6S288, AM10GA, D6S89 and SB1) revealed that there are four haplotypes segregating with disease among the five families analyzed.

4. The trinucleotide repeat is transcribed

To test whether the CAG repeat lies within a gene, reverse transcription-PCR (RT-PCR) was performed using primers immediately flanking the repeat (Rep1 and Rep2) as well as primers which amplify a sequence immediately adjacent to the repeat (Pre1 and Pre2). The RT-PCR analysis confirms that the CAG repeat is present in mRNA from lymphoblasts. Furthermore, northern blot analysis of human poly(A)⁺RNA from various tissues, using a 1.1 kb subclone (C208-1.1) from the 3.36-kb *Eco*RI fragment as a probe, identified a 10 kb transcript which is expressed in brain, skeletal muscle, placenta and to a lesser extent in kidney, lung and heart. The expression of this transcript is considerable in skeletal muscle. When the 3.36-kb *Eco*RI fragment was used as a probe on the northern blot the same size transcript was detected.

5. PCR analysis of the CAG repeat

To confirm that the CAG repeats were involved in the observed length variation, we analyzed the size of PCR-amplified fragments in 45 unaffected spouses and 31 SCA1 affected individuals using synthetic oligonucleotides that flank the CAG repeat. One pair of primers (CAG-a/CAG-b) was located within 9-bp of the repeats and identified length variation indicating that the CAG repeats are the basis of the variation.

Normal individuals displayed 11 alleles ranging from 25 to 36 repeat units (Table 8). Heterozygosity in normal individuals was 84%. Examination of this sequence in 31 individuals affected with SCA1 demonstrated that each was a heterozygote with one allele within the size range seen in the normal individuals and a second expanded allele within a range of 43 to 81 repeat units (Figure 11). Late onset SCA1 individuals showed at least 43 repeats, while 59-81 units were found in the juvenile cases. Figure 12 depicts correlation between the age-at-onset and the number of the repeat units. A linear correlation coefficient (r) of -0.845 was obtained indicating that 71.4% (r^2) of the variation in the age-at-onset can be accounted for by the number of (CAG)_n repeat units. The largest trinucleotide repeat expansion was noted in SCA1 patients with juvenile onset who typically had a more rapid course. It is of interest that all of these patients were offspring of affected males, which is reminiscent of Huntington disease where there is preponderance of male transmission in juvenile cases.

-67-

Sequence analysis of the fragment containing the CAG repeat indicated that there are several extended open reading frames. Translation of the repeat in one of these frames (389-bp) would encode polyglutamine.

5

Table 8.
Comparison of the number of CAG repeat units
on normal and SCA1 chromosomes

10	Number of Repeats	Normal Chromosomes		SCA1 Chromosomes	
		Number	Frequency	Number	Frequency
	≥ 60	0	0	4	0.13
	50 - 59	0	0	17	0.55
15	43 - 49	0	0	10	0.32
	37 - 42	0	0	0	0
	35 - 36	1	0.01	0	0
	30 - 34	49	0.55	0	0
	≤ 29	40	0.44	0	0
20	TOTAL	90	1.00	31	1.00

25 **IV. Isolation of SCA1 cDNA**

A. Methods

1. Screening of cDNA libraries.

Three cDNA libraries were screened: a human fetal brain library from Stratagene (La Jolla, CA), a human fetal brain library constructed in λ -Zap II with the inserts cloned into the *Not*I restriction site (provided by Dr. Cheng Chi Lee at Baylor College of Medicine), and an adult cerebellar cDNA library from Clontech (Palo Alto, CA). The libraries were plated on 150 cm plates at a density

30

of 50,000 pfu per plate using bacterial strain LE392 (ATCC number 33572). Hybond-N filters (Amersham, Arlington Heights, IL) were used to carry out plaque lifts. The fragments used as probes in the first screening included a mixture of two polymerase chain reaction (PCR) products obtained by using the primers Rep1 and Rep2 (Figure 3) immediately flanking the repeat and the primers Pre1 and Pre2 (Figure 3) which amplify a sequence immediately adjacent to the repeat, and a 1.1 kb subclone of the 3.36-kb *EcoRI* fragment (Figure 1). The 1.1 kb fragment (C208-1.1) is located 540 bp 3' to the CAG repeat. A 9-kb *EcoRI* genomic fragment derived from the same cosmids containing the CAG repeat was also used in this screening. Subsequent rounds of screening were carried out on the same libraries using as probes cDNA clones 31-5, 3J, 3c7-2 and 3c7 (Figure 13). Genomic and cDNA probes were labeled using the random priming technique described in A.P. Feinberg et al., Anal. Biochem., **137**, 266-267 (1984). Repetitive sequences were blocked as described in P.G. Sealy et al., Nucl. Acids Res., **13**, 1905-1922 (1985). Briefly, the probes were reassociated with a large excess of shear human placental DNA. The nonrepetitive regions remained single-stranded and no separation of the single-stranded fragments from the reassociated fragments was necessary in order to allow the signal from low copy number components to be detected in subsequent transfer hybridizations. Hybridization of the filters was then carried out following standard protocols as described in H.Y. Zoghbi, et al., Am. J. Hum. Genet., **42**, 877-883 (1988).

2. DNA sequencing and sequence analysis.

Shotgun libraries were constructed in M13 as described in A.T. Bankier, et al., Meth. Enzymol., **155**, 55-93 (1987) for each of the following cDNA clones: 8-8, 31-5, 3c5, 3c7-1, 3J, 3c7-2, 3c7 (Figure 13). Twenty to thirty M13 subclones were sequenced for each cDNA clone using an Applied Biosystem, ABI 370A, automated fluorescent sequencer, as described in R. Gibbs, et al., Proc. Natl. Acad. Sci. U.S.A., **86**, 1919-1923 (1989). Some cDNA clones (8-9b, 8-9a, AX1, B21, B11, 3c28) were partially sequenced manually using a Sequenase sequencing kit (USB, Cleveland, OH) on double-stranded templates, according to the manufacturer's recommendations. The sequence coverage in terms of numbers of cDNA/genomic clones analyzed was 3-4X in the coding and 5'UTR and 2X in the

3'UTR. All RT-PCR, 5'-RACE-PCR and inverse-PCR products were sequenced manually after subcloning into *Sma*I-digested pBluescript SK- plasmid (Stratagene, La Jolla, CA) modified using the T-vector protocol as described in D. Marchuk et al., Nucl. Acids Res., **19**, 1154 (1990). Use of this protocol facilitates cloning. Briefly, *Taq* polymerase ordinarily causes a template-independent addition of adenosine at the 3' end of the PCR product, making blunt end ligations difficult. In the T-vector protocol, a thymidine is added to the 3' end of a digested plasmid. The result is a one-base sticky end complementary to the 3' adenosine in the PCR product, which greatly increases cloning efficiency.

Data base searches were carried out using the GCG software package (Genetics Computer Group, Madison, WI) and the BLAST network service from the National Center for Biotechnology Information (S.F. Altschul, et al., J. Mol. Biol., **215**, 403-410 (1990)). The sequence of the SCA1 transcript has been deposited in Genbank, accession number X79204.

3. Northern blot, RT-PCR and genomic PCR analyses.

The northern blot of poly-(A)⁺ RNA from various human tissues and the poly-(A)⁺ RNA from adult human cerebellum were purchased from Clontech (Palo Alto, CA). Poly-(A)⁺ RNA from human lymphoblastoid cells was prepared by first extracting total RNA using guanidinium thiocyanate, followed by centrifugation in a cesium chloride gradient (P. Chomczynski et al., Anal. Biochem., **162**, 156-159 (1987)). Poly-(A)⁺ RNA was selected using Dynabeads oligo (dT)₂₅ from Dynal (Great Neck, NY). First strand randomly primed cDNA synthesis was carried out using MMLV (murine maloney leukemia virus) reverse transcriptase (BRL, Gaithersburg, MD). This was conducted in a 20 µl reaction mixture containing 3 µg RNA, first strand buffer (50 mM Tris-HCl, pH 8.3, 75 mM KCl, 3 mM Mg Cl₂), (BRL, Gaithersburg, MD), 10 mM dithiothreitol (BRL, Gaithersburg, MD), 1 µM 3' end primer, 0.5 units RNasin (Promega, Madison, WI), 5.0 units MMLV reverse transcriptase (BRL, Gaithersburg, MD), 250 µM each deoxynucleotide triphosphate: dGTP, dATP, dCTP, dTTP. The mixture was incubated for 20 minutes at 37°C then put on ice. A 10 µl aliquot was used for the PCR reaction. First strand randomly primed cDNA from human brain, liver and adrenal were provided by Dr. G. Borsani (Baylor College of Medicine).

-70-

RT-PCR for detection of alternative splicing was carried out with primers 9b and 5R and with primers 5F and 5R (Figure 15) under the following conditions: an initial denaturation step at 94°C for 5' followed by 30 cycles of 94°C for 1 minute, 60°C for 1 minute and 72°C for 2 minutes. The reaction mixture
5 contained 10 µl cDNA, PCR buffer (50 mM KCL, 10 mM Tris-HCl, pH 8.3, 1.25 mM MgCl₂), 1 µM of the relevant 3' primer (primer 5R), 2% formamide and 1.25 units Amplitaq enzyme (Perkin Elmer, Norwalk, CT).

RT-PCR on lymphoblastoid cell lines with primers Rep1 and Rep2 for detection of expression of SCA1 mRNA was carried out using "hot start" PCR
10 with three cycles of: 97°C for 1 minute, 57°C for 1 minute and 72°C for 1 minute. Following that 33 cycles of 94°C for 1 minute, 55°C for 1 minute and 72°C for 1 minute were carried out. Twenty microliters of the PCR reactions was then resolved on a 2% agarose gel (2 g Ultrapure agarose (BRL, Gaithersburg, MD) in 40 mM Tris-acetate, 1 mM EDTA, pH 8.0) and blotted onto Sureblot membrane (Oncor,
15 Gaithersburg, MD). The filter was hybridized with a (GCT)₇ oligonucleotide end-labeled with γ-³²P-ATP. Hybridizations were done in a solution of 1 M NaCl, 1% sodium dodecyl sulfate (SDS) (Sigma Chemical Company, St. Louis, MO) and 10% (w/v) dextran sulphate (Sigma Chemical Company, St. Louis, MO). Filters were washed in 2 x SSC (1 x SSC is 0.15 M sodium chloride and 0.015 M sodium
20 citrate), and 0.1% SDS at room temperature for 15 minutes, followed by a 15 minute wash at room temperature in a solution prewarmed to 67°C.

B. Results

Two human fetal brain cDNA libraries were screened using as probes
25 various DNA fragments from the cosmid clone shown to contain the CAG repeat. Five cDNA clones were identified; these included clone 31-5 containing the CAG repeat, and clone 3J which was found not to overlap with 31-5 (Figure 13). Northern blot analysis revealed that clones 31-5 and 3J identified the same 11-kb transcript detectable in all tissues examined (Figure 14). Accordingly, the same two
30 human fetal brain cDNA libraries and a human adult cerebellar cDNA library were used for several rounds of screening in order to obtain the full length transcript. As a result, 22 cDNA clones were isolated and characterized by sequence and PCR analyses to assemble a contig spanning the SCA1 transcript. Twelve of the phage

clones spanning the cDNA contig are shown in Figure 13. These clones were sequenced allowing the assembly of the entire sequence of the SCA1 cDNA which spans 10,660 bp (Figure 15).

Sequence analysis revealed a coding region of 2448 bp starting with a putative ATG initiator codon at base 936 located within a nucleotide sequence that fulfills Kozak's criteria for an initiation codon (M. Kozak, *J. Cell. Biol.*, **108**, 229-241 (1989)). An in-frame stop codon is present 57 bp upstream of that ATG in three independent cDNA clones as well as in genomic DNA. Furthermore, both the ATG at the beginning of the coding region and the upstream stop codon have been found in the murine homologue of SCA1 in the murine fetal brain library (Stratagene, La Jolla, CA). The SCA1 gene therefore encodes a polypeptide of about 816 amino acids, with an expected size of 87 kD, designated ataxin-1. However, one cannot exclude the possibility that the coding region begins at any of the other ATGs, located downstream of the first methionine, which would result in a smaller protein.

The CAG repeat is located within the coding region 588 bp from the first methionine and encodes a polyglutamine tract. The open reading frame ends with a TAG stop codon at base 3384. Therefore, this transcript has a 5' untranslated region (5'UTR) of 935 bp and a 3' untranslated region (3'UTR) of 7277 bp. The transcript ends with a tail of 57 adenosine residues; a polyadenylation signal, AATAAA, is found 23 nucleotides upstream of the poly(A) tail. Homology searches using both the DNA sequence of the coding region and the predicted protein sequence (lacking the CAG repeat and the polyglutamine tract, respectively) revealed no significant homology with other known proteins in the data base. Analysis of the sequence of ataxin-1 failed to reveal the presence of any strong phosphorylation sites as well as any specific motifs such as DNA binding or RNA binding domains. The putative secondary structure of this protein is compatible with that of a soluble protein as no hydrophobic domains were identified. A DNA sequence data base search revealed an identity between 380 bp in the 3'UTR of the SCA1 transcript and an expressed sequence tag (EST04379) isolated from a human fetal brain cDNA library (M.D. Adams, M.D. et al., *Nature Genet.*, **4**, 256-267 (1993)).

V. Organization of the SCA1 Transcript: Evidence for Alternative Splicing in the 5'UTR

A. Methods

1. 5'-RACE-PCR

5 First strand cDNA was prepared from 1 mg of poly-(A)⁺ RNA from human adult cerebellum (Clontech, Palo Alto, CA) using the primer 5R (Figure 15) as described in Example IV. 5'-RACE-PCR was carried out as described in M.A. Frohman in PCR Protocols. A Guide to Methods and Applications; M.A. Innis, et al., Eds.; Academic Press: San Diego (1990) using SCA1 primers 5a and X4-1
10 (Table 9) as specific primers. The product was then electrophoresed through a 1.2% agarose gel, blotted onto SureBlot hybridization membrane (Oncor, Gaithersburg, MD) as described in Example II above, and then, to test the specificity of the product, hybridized to a SCA1 specific probe represented by a PCR product spanning 118 bp between primer 9b in exon 1 and primer X3-1 (Table 9) in exon 3.

15

Table 9.

Primer sequences for inverse-PCR

Exon	Primer 1	Primer 2
2	X2-1 (181-164) GTAGTAGTTTTGTGAGG	X2-2 (185-203) CACCAAGCTCCCTGATGGA
3	X3-1 (246-229) GCTTGAATGGACCACCCT	X3-2 (277-296) ATCTCCTCCTCCACTGCCAC
4	X4-1 (347-329) AGACTCTTTCACATGCTC	X4-4 (407-425) TTCAGCCTGCACGGATGGT
5	5a (482-463) TGGCAGTGGAGAATCTCAGT	5-2 (519-538) TGCTGCAAGGAACTGATAGC
6	10a (598-580) AATGGTCTAATTCTTTGG	10b (607-625) GAGAAAGAAATCGACGTGC
7	6-1 (714-695) ACAGGCTCTGGAGGGCTCCT	X5-2 (723-742) TCCATGGTGAAGTATAGGCT
9	9-1 (2919-2900) AGCAGGATGACCAGCCCTGT	9-2 (2939-2957) GCTCTTTGATTTGCCGTGT
All primers are read in the 5' to the 3' direction. Numbers in parenthesis represent the coordinates of each primer within the SCA1 cDNA sequence (Figure 15).		

B. Results

To characterize the genomic region flanking the CAG repeat, the 3.36-kb *Eco*RI genomic fragment known to contain this repeat was completely sequenced. Alignment of this genomic sequence with the cDNA sequence allowed us to determine that the 3.36-kb *Eco*RI fragment contains a 2080-bp exon which has 160 bp of 5'UTR, the first potential initiation codon and the first 1920 bp of the coding region. The rest of the coding region lies within the next downstream exon as detected by PCR analysis on genomic DNA. The last coding exon, which maps to a 9-kb *Eco*RI fragment in genomic DNA also contains 7277 bp of 3'UTR for a total length of 7805 bp (Figure 16a).

Evidence for alternative splicing in the 5'UTR was initially suggested based on the hybridization pattern of the two most 5' cDNA clones, 8-8 and 8-9b (Figure 13) to Southern blots containing *Eco*RI-digested genomic DNA from total human DNA and YACs spanning the SCA1 region. At least three strongly hybridizing fragments in addition to the 3.36-kb *Eco*RI fragment were seen. As neither of the cDNA clones contains an *Eco*RI site, this result suggested the presence of several exons in the 5'UTR of the SCA1 transcript. Given these data and the unusual length of the 5'UTR, this region was characterized in more detail.

Alignment analysis of the sequence of clones 8-8 and 8-9b revealed the presence of two different 5' sequences diverging at basepair 322. This result was highly suggestive of alternative splicing. In order to test this hypothesis, reverse transcription-PCR (RT-PCR) was performed on mRNA from cerebellar tissue using the primers indicated in Figure 15. When the primers 9b (specific for 8-9b clone) and 5R (present in both clones) were used in the RT-PCR analysis three products were obtained: one of the expected size (246 bp) and at least two fragments of larger size (Figure 16b). The same result was obtained when RT-PCR was carried out on liver, adrenal, brain and lymphoblast cDNAs. The various RT-PCR products were cloned and sequenced. Sequence analysis of all these products and comparison with the sequence of phage clones 8-8 and 8-9b confirmed that they were the result of alternative splicing. Figure 16a shows the structure of all the cDNA clones which contain the 5' exons of the SCA1 gene and depicts the splice variants. Based on sequence analysis of three cDNA clones and characterization of cerebellar RT-PCR products, five exons (exons 1 through 5) were identified and

their borders in the transcript were determined. Exons 2, 3 and 4 are alternatively spliced in the clones examined and in cerebellar tissue, whereas exon 5 was present in all the cDNA clones and RT-PCR products.

Rescreening of cDNA libraries with clones 8-8 and 8-9b as probes did not yield any additional cDNA clones. To identify additional alternatively spliced exons in the 5'UTR and to confirm initial results, 5'-RACE-PCR was carried out on reverse transcribed cerebellar mRNA using primers from the 5' end of exons 5 and 4. A 218-bp product was identified and its specificity was confirmed by Southern analysis using an internal PCR product as probe. Sequence analysis of the 5'-RACE-PCR product, furthermore, confirmed the alternative splicing of two exons (2 and 3) and allowed the identification of an additional 127 bp at the 5' end of this gene (Figure 16a).

VI. Identification of Intron-Exon Boundaries and Determination of the Genomic Structure of SCA1

A. Methods

1. Identification of intron-exon boundaries

The boundaries of exons 2-9 were identified by inverse-PCR. To carry out inverse-PCR, YAC agarose plugs were digested to completion as described in M.C. Wapenaar, et al., Hum. Mol. Genet., 2, 947-952 (1993) using frequent-cutter restriction enzymes such as *Sau3aI*, *TaqI*, *HaeIII* and *MspI* purchased from Boehringer Mannheim Biochemicals (Indianapolis, IN) and used as recommended by the manufacturer. The plugs were then digested with β agarase I (USB, Cleveland, OH) following the manufacturer's recommendations and subsequently phenol-chloroform (Boehringer Mannheim Biochemicals, Indianapolis, IN) extracted, precipitated with ethanol and resuspended in 12 ml of TE (TE: 10 mM Tris-HCl, 1 mM EDTA) pH 8. Fifty ng of DNA from each digest was then circularized according to the published protocol of J. Groden et al., Cell, 66, 589-600 (1991). Diverging PCR primers were designed within the cDNA and used on the circularized product under the amplification conditions described in J. Groden et al., Cell, 66, 589-600 (1991). PCR products were then subcloned and sequenced as described in Example II, above. Inverse-PCR identified all intron/exon boundaries except the boundary of exon 1. Accordingly, a 9-kb *EcoRI*

genomic fragment found to contain exon 1 was subcloned from a cosmid derived from YAC 227B1. (Example II). This subclone was subsequently partially sequenced to identify the boundary of exon 1.

5 2. Mapping of cDNA clones to the YACs and cosmids.

Southern blots containing *Eco*RI-digested DNAs from YACs spanning the SCA1 critical region as well as Southern blots containing DNAs from the YACs digested with rare-cutter enzymes (see previous section) were hybridized, using the standard protocol described in H.Y. Zoghbi et al., *Am. J. Hum. Genet.*, **42**, 877-883 (1988), to various SCA1 cDNA clones and to all the genomic fragments containing the intron-exon boundaries. Briefly, restriction fragments were separated by electrophoresis on 0.7% agarose gels, denatured and transferred to Nytran (Schliecher and Schuell, Keene, NH) filters. Probes were ³²P-labeled using the oligohexamer labeling method (A.P. Feinberg et al, *Anal. Biochem.*, **132**, 6-13 (1983)). After hybridization the filters were washed and autoradiography was performed, as described in Zoghbi et al., *Am. J. Hum. Genet.*, **42**, 877-883 (1988).

B. Results

Complete sequencing of the 3.36-kb *Eco*RI fragment provided the intron-exon boundaries for the 2080-bp exon containing most of the coding region (Figure 17). In order to determine the actual number of exons and to obtain all of the intron-exon boundaries, an inverse-PCR strategy was adopted using two overlapping YAC clones, 227B1 and 149H3, known not to contain any rearrangements (see Example II). A total of nine exons, seven of which are in the 5'UTR, were identified and splice junctions for exons 1 through 9 were subcloned and sequenced (Figure 17). The schematic on top of Figure 16a shows the nine exons and their respective sizes. In the 5' untranslated region, alternative splicing involves exons 2, 3 and 4, but not exons 5, 6 and 7 in over 5 phage cDNA clones analyzed. The putative exon 1 encompasses 157 bp and hybridizes very strongly to an *Eco*RI fragment derived from hamster genomic DNA.

To study the genomic organization of the SCA1 gene, ten cDNA clones as well as genomic fragments containing the splice junctions for all the exons were mapped by Southern analysis and localized on a long range restriction map of

four overlapping YAC clones spanning the SCA1 critical region (Figure 18). This analysis revealed that the gene spans at least 450 kb of genomic DNA and that the putative first exon maps to a genomic fragment containing a hypomethylated CpG island. Detailed restriction analysis of the intron between the two coding exons (8 and 9) revealed that this intron is approximately 4.5-kb in length. The sizes of the remaining introns were estimated from the long range restriction map and by PCR analysis and ranged from 650 bp (intron 2) to nearly 200 kb (intron 7) (Figure 18).

VII. Expression of the SCA1 mRNA in SCA1 Patients

As a first step toward understanding the mechanism by which the expansion of a trinucleotide CAG repeat leads to neurodegeneration in SCA1, the level of transcription of SCA1 from the expanded alleles in patients was investigated. RT-PCR was carried out with primers Rep1 and Rep 2 which flank the CAG repeat as described in Example V using lymphoblastoid mRNAs from SCA1 patients with repeat sizes ranging from 43 to 69. This analysis revealed that mRNA was expressed from both the normal allele and the expanded allele (Figure 19).

VIII. Cloning of portions of the SCA1 Gene into the pMALTM-2 Vector

DNA from the SCA1 gene was cloned into the pMALTM-c2 vector (New England Biolabs, Beverly, MA), which produces a chimeric protein consisting the maltose-binding protein fused to the N-terminus of the protein of interest (ataxin-1) in a linkage that can subsequently be conveniently cleaved. To obtain DNA for cloning, SCA1 DNA was amplified and isolated clone 31-5 (Figure 13) using standard PCR techniques. The manufacturer's instructions were followed in designing the appropriate oligonucleotide primers (pMALTM vector Package Insert, 1992 New England Biolabs, revised 4/7/92). In each case an *EcoRI* linker site was designed into the 5' primer and a *HindIII* linker site was designed into the 3' primer to facilitate cloning. Three different amplification products were obtained. In one, DNA was isolated utilizing two 20-mer PCR primers COD and RCOD (Table 10) that hybridized to the 5' and 3' ends of the coding regions, such that the stretch of DNA being amplified contained residues presumed to encode the entire sequence of ataxin-1, beginning with Met1 and ending with Lys 817 (Figure 15). The amplified product was then cloned into the *EcoRI/HindIII* site in the polylinker region of in

-77-

pMALTM-c2 following instructions provided by the manufacturer. Two other constructs were made in the same way using PCR to isolate shorter segments of DNA. In both cases the same 3' end primer was used, but different 5' primers were employed (Table 10). One 5' primer (3COD) was designed such that the amplified product began at Met277 (the fourth methionine in the coding region), the other 5' primer (8COD) such that the amplified product began at Met548. pMALTM-c2 was transformed into competent cells containing a lacZ Δ M15 allele for α -complementation and cultured as recommended by the manufacturer.

10

Table 10.
Primers for Cloning Into pMal Vector

<u>Primer Name</u>	<u>Nucleotide Sequence</u>
COD	TGT GAA TTC ATG AAA TCC AAC CAA GAG CG
3COD	TGT GAA TTC ATG ATC CCA CAC ACG CTC AC
8COD	TGT GAA TTC ATG GTG CAG GCC CAG ATC
RCOD	TTC GAA GCT TCT ACT TGC CTA CAT TAG AC

15 **IX. Expression of Ataxin-1, Design of Antigenic Peptides and Production of Antibodies**

The fusion protein expressed by the constructs in Example VII were purified as directed by the manufacturer using affinity chromatography (pMALTM vector Package Insert, 1992 New England Biolabs, revised 4/7/92). The purified protein was electrophoresed using 8% SDS polyacrylamide electrophoresis and electroeluted. The best expression (about 27 mg from 1 L of cells) was obtained from the shortest construct, but all constructs produced measurable levels of protein of a size consistent with their respective cloned gene product.

20 Antibody response in rabbits was initiated using the multiple antigenic peptide strategy of V. Mehra et al., Proc. Natl. Acad. Sci. USA, 83, 7013-7017 (1986). In addition to the three electroeluted cloned gene products described in the preceding paragraph, three synthetic peptides were used as well. The synthetic peptides used were Peptide A (amino acids 4 through 18), Peptide B

-78-

(amino acids 162 through 176) and Peptide C (amino acids 774 through 788). These peptides were chosen such that they showed little or no homology with other known short amino acid stretches in proteins and also such that they contained proline, which makes it more likely that these fragments are located on the surface of the protein, thus making it more likely that antibodies to the fragments will react with the whole protein as well.

Immunoglobulin (IgG) from rabbit blood was purified, and antibody/antigen results were analyzed using Western blots as described in Gershoni et al., *Anal. Bioch.*, 131, 1-15 (1983). IgG from rabbits injected with the cloned gene products and the synthetic sequences were found to hybridize to their respective antigens. The anti-sera from rabbits immunized with the 8COD-RCOD gene product (i.e., the ataxin-1 fragment spanning residues 548 through 817) hybridized with a protein of the expected size in brain tissue extracts from mouse, rats, and humans. A similar size protein has also been detected using lymphoblasts. This hybridization is blocked by preincubation with the polypeptide antigen, and not blocked by unrelated antigens. In particular, antibodies raised against Peptide C are blocked by either Peptide C or the short gene product.

X. Molecular and Clinical Correlations in Spinocerebellar ataxia type 1 (SCA1)

A. Materials and Methods

1. Family Material

Members representing 87 kindreds with dominantly inherited ataxia were evaluated. Nine kindreds of diverse ethnic background (Caucasian American, African American, South African, Siberian Iakut) were already known to have SCA1 based on linkage to the HLA locus and to D6S89 on chromosome 6p. Genotypic analysis of the SCA1 CAG repeat was carried out on all nine kindreds to determine if all known SCA1 families had the same mutational mechanism involving repeat expansion. Most of the study participants were personally examined. The affected status was always confirmed by a neurologist, but the age of onset was based on historical information from the patient and/or other family members. Severity of disease was measured by the age at death minus the age of onset. Detailed characterization of the repeat variability was carried out for all nine

-79-

kindreds. To identify additional kindreds with a CAG expansion at the SCA1 locus, affected individuals from 78 newly identified families with dominantly inherited ataxia were clinically examined. Blood was collected from at least one affected individual from each of these kindreds and screened by DNA analysis for the presence of a CAG repeat size within the expanded range (≥ 42 repeats). Although there was no evidence that these 78 individuals are related, there is a chance that some of the affected patients come from the same families.

To assess the distribution of CAG repeat sizes on normal chromosomes further, the number of CAG repeats was determined for 304 normal chromosomes from unrelated individuals of various ethnic backgrounds.

2. Molecular Studies

Blood samples were used to establish lymphoblastoid cell lines by Epstein-Barr virus transformation. Genomic DNA was isolated either directly from venous blood or from lymphoblastoid cell lines. Blood samples were collected from these patients over an 8-year period, during which time 29 patients died. PCR reactions were performed using the Rep1 (TTGACCTTTACACCTGCAT) and Rep2 (CAACATGGGCAGTCTGAG) primers. Fifty nanograms of genomic DNA was mixed with 5 pmol of each primer in a total volume of 20 μ l containing 1.25 mM $MgCl_2$, 250 μ M dNTPs, 50 mM KCl, 2% formamide, 10 mM Tris-HCl pH 8.3 and 1 unit ampliTaq (Perkin-Elmer/Cetus). The Rep1 primer was labelled at the 5' end with [γ - ^{32}P]ATP. Samples were denatured at 94°C for 4 minutes, followed by 30 cycles of denaturation (94°C, 1 minute), annealing (55°C, 1 minute) and extension (72°C, 2 minutes). Six μ l of each PCR reaction was mixed with 4 μ l formamide loading buffer, denatured at 90°C for 2 minutes, and electrophoresed through a 6% polyacrylamide/7.65 M urea DNA sequencing gel. Allele sizes were determined by comparing migration relative to an M13 sequencing ladder.

3. Statistical Analyses

The relationship between age of onset and CAG repeat number on both the affected and the normal chromosomes of patients was evaluated through linear regression analyses. Similarly, the relationship between repeat length and duration of disease was quantified. Ages of onset were used directly in these

analyses, but also following logarithmic and square root transformation. Although the latter transformation provided the best approximation to a normal distribution, results obtained were consistent between analyses before and after transformation. Analysis of variance was performed to detect differences among the families in the mean age of onset, after correction for the effect of the CAG repeat number on age of onset. In addition, the sex of the transmitting parent was included as a possible explanatory variable for variations in age of onset. All regression and variance analyses were carried out with the SPSS package of computer programs, versions 4.0.1.

B. Results

1. Family Studies

All affected individuals from the nine known SCA1 kindreds had an expanded trinucleotide repeat on one of their alleles. No repeat expansions were observed among eight kindreds previously shown by linkage analyses not to be SCA1. These eight kindreds were examined for the SCA1 gene expansion to confirm the linkage results.

Among the 70 other dominant ataxia families analyzed, three (4%) were found to have an expanded CAG repeat on one of the SCA1 alleles. Of all of the dominant kindreds studied, 12 of 87 (14%) have an expanded CAG repeat at the SCA1 locus. While the sample size is relatively small, and both estimates are arguably biased to exclude or select for SCA1 kindreds, expanded CAG repeat tracts within the SCA1 gene clearly account for only a small fraction of this complex group of diseases. The distribution of the CAG repeat number from normal controls and from ataxic individuals that did not have an expansion were similar (data not shown). These data argue against the involvement of the CAG repeat at the SCA1 locus in these families. However, it is still possible that some of these small families have other mutations at the SCA1 locus.

The typical clinical findings in the genetically proven SCA1 kindreds were gait and limb ataxia, dysarthria, pyramidal tract signs (spasticity, hyperreflexia, extensor plantar responses) and variable degrees of oculomotor findings which include one or more of the following: nystagmus, slow saccades, and ophthalmoparesis. In the later stages of the disease course, bulbar findings consistent

-81-

with dysfunction of cranial nerves IX, X, and XII became evident. Also, dystonic posturing and involuntary movements including choreoathetosis became apparent in the later stages of the disease. Motor weakness, amyotrophy, and mild sensory deficits manifested as proprioceptive loss were also detected. Although ataxia, dysarthria and cranial nerve dysfunction were consistently present in every SCA1 affected individual, considerable intrafamilial variability was noted with regard to all of the other clinical features. Juvenile onset (≤ 18 years) was observed in four kindreds. Of interest is the finding that juvenile onset cases typically inherited the disease gene from an affected father. Several of the kindreds that did not have an expanded SCA1 CAG repeat, displayed the same clinical findings as those observed in SCA1 kindreds confirming the inherent difficulty in clinically classifying this group of disorders. While it is possible that some of these kindreds have other mutations at the SCA1 locus, the disease locus (loci) for eight of these families has also been excluded from the SCA1 region by linkage analyses.

2. Repeat Analysis on Normal and SCA1 Chromosomes

Figure 20 shows the size distribution of the CAG repeats on 304 chromosomes from unaffected control individuals who are at risk for ataxia, and 113 expanded alleles from individuals affected with the disease. The normal alleles range in size from 19 to 36 CAG repeat units. Over 95% of the normal alleles contain from 25 to 33 CAG repeat units, the majority (65%) of which contain 28 to 30 repeats. The mean repeat size on normal chromosomes for the African Americans, Caucasian, and South African populations are very similar with 29.1, 29.8, and 29.4 CAG repeat units, respectively. Combined heterozygosity for the CAG repeat at the SCA1 locus was 0.809 for the populations examined, giving an overall polymorphism information content (P.I.C.) value of 0.787. No change in CAG repeat length was observed for 135 meioses of SCA1 alleles containing CAG repeat tracts within the normal range, i.e., all were inherited in a Mendelian fashion. In contrast, 41 of the 62 meioses involving expanded SCA1 alleles changed in repeat size. The rate of repeat instability for female meioses is 60% while the instability observed for males was 82%.

-82-

The number of CAG repeats found on SCA1 chromosomes from 113 affected individuals was always greater than the number of repeats on normal chromosomes, ranging from 42 to 81 with a means of 52.6 (Figure 20).

5 All patents, patent documents, and publications cited herein are incorporated by reference. The foregoing detailed description and examples have been given for clarity of understanding only. No unnecessary limitations are to be understood therefrom. The invention is not limited to the exact details shown and described, for variations obvious to one skilled in the art will be included within the
10 invention defined by the claims.

WHAT IS CLAIMED IS:

1. A nucleic acid molecule containing a CAG repeat region of an isolated autosomal dominant spinocerebellar ataxia type 1 (SCA1) gene, said gene located within the short arm of chromosome 6.
2. The nucleic acid molecule of claim 1 corresponding to the entire SCA1 gene.
3. The nucleic acid molecule of claim 1 wherein the SCA1 gene encodes ataxin-1.
4. The nucleic acid molecule of claim 3 of about 2.4-11 kb in length containing the coding region of the SCA1 gene.
5. The nucleic acid molecule of claim 1 wherein the CAG repeat region is represented by $(CAG)_n$ and $n = 2-36$.
6. The nucleic acid molecule of claim 5 wherein $n = 19-36$.
7. The nucleic acid molecule of claim 1 wherein the CAG repeat region is represented by $(CAG)_n$ and $n > 36$.
8. The nucleic acid molecule of claim 7 wherein $n \geq 43$.
9. The nucleic acid molecule of claim 1 wherein the molecule is a single-stranded polynucleotide.
10. The nucleic acid molecule of claim 9 wherein the single stranded polynucleotide is cDNA.
11. The nucleic acid molecule of claim 9 wherein the single stranded polynucleotide is mRNA.

12. The nucleic acid molecule of claim 1 wherein the nucleic acid is genomic DNA.
13. An isolated oligonucleotide that hybridizes to a nucleic acid molecule containing a CAG repeat region of an isolated SCA1 gene; said oligonucleotide having at least about 11 nucleotides.
14. The isolated oligonucleotide of claim 13 having at least about 16 nucleotides.
15. The isolated oligonucleotide of claim 14 having no more than about 35 nucleotides.
16. The isolated oligonucleotide of claim 13 that produces a primed product of about 70-350 base pairs.
17. The isolated oligonucleotide of claim 16 that produces a primed product of about 100-300 base pairs.
18. The isolated oligonucleotide of claim 13 that hybridizes to the nucleic acid molecule within about 150 nucleotides on either side of the CAG repeat region.
19. The isolated oligonucleotide of claim 18 that hybridizes to the nucleic acid molecule directly adjacent to the (CAG)_n region.
20. The isolated oligonucleotide of claim 13 having at least about 100 nucleotides.
21. The isolated oligonucleotide of claim 20 having at least about 200 nucleotides.
22. The isolated oligonucleotide of claim 13 comprising a nucleotide sequence selected from the group consisting of CCGGAGCCCTGCTGAGGT (CAG-a), CCAGACGCCGGGACAC (CAG-b), AACTGGAAATGTGGACGTAC (Rep-1), CAACATGGGCAGTCTGAG (Rep-2),

-85-

CCACCACTCCATCCCAGC (GCT-435), TGCTGGGCTGGTGGGGGG
(GCT-214), CTCTCGGCTTTCTTGGTG (Pre-1), and
GTACGTCCACATTTCAGTT (Pre-2).

23. A method for detecting the presence of a DNA molecule containing a CAG repeat region of the SCA1 gene comprising:
 - (a) digesting genomic DNA with a restriction endonuclease to obtain DNA fragments;
 - (b) probing said DNA fragments under hybridizing conditions with a detectably labeled gene probe, which hybridizes to a nucleic acid molecule containing a CAG repeat region of an isolated SCA1 gene having at least about 11 nucleotides;
 - (c) detecting probe DNA which has hybridized to said DNA fragments; and
 - (d) analyzing the DNA fragments for a CAG repeat region characteristic of the normal or affected forms of the SCA1 gene.
24. The method of claim 23 wherein the step of analyzing comprises analyzing for a (CAG)_n region wherein $n > 36$.
25. The method of claim 24 wherein the step of analyzing comprises analyzing for a (CAG)_n region wherein $n \geq 43$.
26. The method of claim 23 wherein the detectably labelled DNA sequence comprises a portion of an *Eco*RI fragment of the SCA1 gene.
27. The method claim 26 wherein the *Eco*RI fragment comprises about 3360 base pairs.

28. A method for detecting the presence of a DNA molecule located within an affected allele of the SCA1 gene comprising:
 - (a) treating separate complementary strands of a DNA molecule containing a CAG repeat region of the SCA1 gene with a molar excess of two oligonucleotide primers;
 - (b) extending the primers to form complementary primer extension products which act as templates for synthesizing the desired molecule containing the CAG repeat region;
 - (c) detecting the molecule so amplified; and
 - (d) analyzing the amplified molecule for a CAG repeat region characteristic of the SCA1 disorder.
29. The method of claim 28 wherein the step of analyzing comprises analyzing for a $(CAG)_n$ region wherein $n > 36$.
30. The method of claim 29 wherein the step of analyzing comprises analyzing for a $(CAG)_n$ region wherein $n \geq 43$.
31. A protein encoded by the SCA1 gene having therein a glutamine repeat region.
32. The protein of claim 31 having a molecular weight of about 20-90 kD.
33. The protein of claim 31 having the amino acid sequence shown in Figure 15.
34. An antibody to a protein encoded by DNA containing a CAG repeat region of the SCA1 gene.
35. A method for detecting the SCA1 disorder comprising:
 - (a) contacting an antibody to a protein encoded by the SCA1 gene with a biological sample containing antigenic protein to form an antibody-antigen complex;
 - (b) isolating the antibody-antigen complex; and

-87-

- (c) sequencing the antigen portion of the antibody-antigen complex using amino acid sequencing techniques.

1/23

FIGURE 1

```

1 TTTTGAAACT TGCAGAGAAC AGGATTATTT CTGGCGGCCT CTGCTGAGTT GCGGTGTGTG
61 TGTGTGTTTG TGTGTGTGTG TATTAGGGAG AGGAAATCGT AGGTCCAGTG TGGACCCAGA
121 GCTAAGGGGA ATCTTGGAGA GTAGTGGCTC TGGCAGATGA GGATTCAGAA ATCGAGTGCA
181 AGGACTGTTT TGGACTTTCA CTGCTAACCT GCTTTTCTC AGTGCCCTGGC TCTGAGGGCA
241 GGGTCCAGCT GGTGTCATGC TCTCCAAGCG CTTCATTTTA TGTTCAGCC AGGCAAAGGA
301 GAGGTGAGAA ATGGAACCAA CATTTCTGAA AAGGAAATTT AAGAACTGCA TCATCTGCCC
361 TTGAAGAAGA AAAGGAGAAA AAAAAACAGG AGAGAGGGTA TTGAGAACAT CTTAGGGGAG
421 TTGTTAACTC CATTAAAAAA TATATGTGTT ACAGTGTTC CTTGCCCAGT GTCTTCATAA
481 TCTTCCTTTA TAATGTGCAG CTGCCACGGC TAGTGTTTTT GTTTTGTGTT TTGTTGTTTT
541 GTTTCGTTTT TGGAGACAGA GTGTCGCTCT GTTGCCCAGG CTGGAGTACA ATGGTGCAT
601 CTCGGCTCAC TGCAACCTCT GCCTCCTGGG TCTCAAGCAAT TCTCCTGCCT CAGCCTCTCA
661 AGTAGCTGGG ACTACAGCCG TGTGCCAGCT AATGTTACAC CAGGCTAAAT TTGTTTTTTA
721 TTTTTTATTT TTGGTAGAGA CGGGGTTTCA CCATGTTAGC CAGGATGGTC TTAATCTCCT
781 GACCTCGTGA TCTGCCTGCC TCGGCCTCCC AAAGTGTGG CTAGTGTTTT CTCTGCTTCA
841 GTGCTTGGGG TATGATTGGG TTATGGGAGT TCACACCGAG TCCAGGGCCT AGTCTTAATC
901 TTGCCAAAGA TGTTCTTTCC CCGGTGCTCA TGTCTGATG TCCTTTCCCT CCTTCCCTTT
961 CTCCTCCCTT TCCTTTTCCC TTTGTCACAT CCCTCTTCCC TTTCCCAGCA TCCAGAGCTG
1021 CTGTTGGCGG ATTGTACCCA CGGGGAGATG ATTCCTCATG AAGAGCCTGG ATCCCTACA
1081 GAAATCAAAT GTGACTTTCC GTTTATCAGA CTAAATCAG AGCCATCCAG AACAGTGAA
1141 CAGTCACCGT GGAGGGGGGA CGGCGAAAAA TGAAATCCAA CCAAGAGCGG AGCAACGAAT
1201 GCCTGCCTCC CAAGAAGCGC GAGATCCCCG CCACCAGCCG GTCCTCGGAG GAGAAGGCC
1261 CTACCTTGAC CCAGCGACAA CCACCGGGTG GAGGGCACAG CATTGGCTCC CGGGCAACCC
1321 TGGTGGCCGG GCGGAGGCA TGGGCCGGCA GGGACCTCGG TGGAGCTTGG
1381 TTTACAACAG GGAATAGGTT TACACAAAGC ATTGTCCACA GGGCTGGACT ACTCCCCGCC
1441 CAGCGCTCCC AGGTCTGTCC CCGTGGCCAC CACGCTGCCT GCCGCGTAGC CCACCCCGCA
1501 GCCAGGGACC CCGGTGTCCC CCGTGCAGTA CGCTCACCTG CCGCACACCT TCCAGTTCAT
1561 TGGGTCTCTC CAATACAGTG GAACCTATGC CAGCTTCATC CCATCACAGC TGATCCCCC
1621 AACCGCCAAC CCCGTCACCA GTGCACTGGC CTCGGCGCAG GGGCCACCAC TCCATCCAG
1681 CGCTCCAGC TGGAGGCCTA TTCCACTCTG CTGGCCAACA TGGGCAGTCT GAGCCAGACG
1741 CCGGGACACA AGGCTGAGCA GCAGCAGCAG CAGCAGCAGC AGCAGCAGCA CCTCAGCAGG
1801 CATCAGCAGC AGCAGCAGCA GCAGCAGCAG CAGCAGCAGC AGCAGCAGCA CCTCAGCAGG
1861 GCTCCGGGGC TCATCACCCC GGGTCCCCC CAACCAGCCC AGCAGAACCA GTACGTCCAC
1921 ATTTCCAGTT CTCCGCAGAA CACCGGCCGC ACCGCCTCTC CTCCGGCCAT CCCCCTCCAC
1981 CTCCACCCCC ACCAGACGAT GATCCCACAC ATGCTCACCC TGGGGCCCCC CTCCCAGGTC
2041 GTCATGCAAT ACGCGACTC CGGCAGCCAC TTTGTCCCTC GGGAGGCCAC CAAGAAAGCC
2101 GAGAGCAGCC GGCTGCAGCA GGCCATCCAG GCCAAGGAGG TCCTGAACGG TGAGTTGGAG
2161 AAGAGCCGGC GGTACGGGGC CCCGTCTCTA GCCGACCTGG GCCTGGGCAA GGCAGGCGGC
2221 AAGTCGGTTC CTCACCCGTA CGAGTCCAGG CACGTGGTGG TCCACCCGAG CCCCTCAGAC
2281 TACAGCAGTC GTGATCCTTC GGGGGTCCGG GCCTCTGTGA TGGTCTGTGC CAACAGCAAC
2341 ACGCCCGCAG CTGACCTGGA GTGCAACAG GCCACTCATC GTGAAGCCTC CCCTTCTACC
2401 CTCAACGACA AAAGTGGCCT GCATTTAGGG AAGCCTGGCC ACCGGTCTTA CGCGCTCTCA
2461 CCCCACACGG TCATTAGAC CACACACAGT GCTTCAGAGC CACTCCCGGT GGACTGCCAG
2521 CCACGGCCTT CTACGCAGGG ACTCAACCCC CTGTATCGG CTACCTGAGC GGCCAGCAGC
2581 AAGCAATCAC CTACGCCGGC AGCCTGCCCC AGCACCTGGT GATCCCCGGC ACACAGCCCC
2641 TGCTCATCCC GGTCCGCAGC ACTGACATGG AAGCGTCGGG GGCAGCCCCG GCTCATAGTCA
2701 CGTCATCCCC CCAGTTTGCT CGAGTGCCTC ACACGTTCTG CACCACCGCC CTTCCTCAAG
2761 GCGAGAACTT CAACCTGAG GCCCTGGTCA CCCAGGCCGC CTACCCAGCC ATGGTGCAGG
2821 CCCAGATCCA CCTGCCTGTG GTGCACTCCG TGGCCTCCCC GCGGGCGGCT CCCCCTACGC
2881 TGCCTCCCTA CTTTCATGAA GGCTCCATCA TCCAGTTGGC CAACGGGGAG CTAAAGAAGG
2941 TGGAAGACTT AAAACAGAAG ATTTATCCA GAGTGCAGAG ATAAGCAAGC ACCTGAAGAT
3001 CGACTCCAGC ACCGTAGAGA GGATTGAGA CAGCCATAGC CCGGGCTGGG CCGTGATACA
3061 GTTCGCCGTC GGGGAGCACC GAGCCAGGT AACGTTAGCC AGGGTGGCAC AGGGATGGGA
3121 CACCATACCG TGATGCCATC ATCATCTCCT GGCAAGACGA ATTGCTTCTA TGAGGCAGGA
3181 TTAAGGGTTC TCGGGTACAC CTAGACCTTA GACTCGGCCT TTCCCAACTG CGTTCTCTAG
3241 AAAAAATAAG CCCCATTTC CCGTGATCTC TGCTGTGTGT AATGAATTAA CCTCCATGCA
3301 TGGAGAGTGG GGCTAGTTAT GGAGTCCTTG AGACAATCCA GAAACTCACC ACTCTCGTTA
3361 TTTTTT

```


3/23

1 GATCCCCCACC ACCGCCAACC CCGTCACCAG TGCAGTGGCC TCGGCGCAGG
 GCT-435 \

51 GGCCACCACT CCATCCCAGC GCTCCCAGCT GGAGGCCTAT TCCACTCTGC
 Rep-2 \ CAG-b \

101 TGGCCAAACAT GGCAGTCTG AGCCAGACGC CGGGACACAA GGCTCAGCAG
 Rep-2 \

151 CAGCAGCAGC AGCAGCAGCA GCAGCAGCAG CAGCATCAGC ATCAGCAGCA
 / CAG-a

201 GCAGCAGCAG CAGCAGCAGC AGCAGCAGCA GCAGCAGCAG CTCAGCAGGG
 / GCT-214 \

251 CTCCGGGGGCT CATCACCCTCG GGTCCCCCCC ACCAGCCCAG CAGAACCAGT
 ← Rep-1 Pre-2 → \

301 ACGTCCACAT TTCCAGTTCT CCGCAGAACA CCGGCCGCAC CGCCTCTCCT

351 CCGGCCATCC CCGTCCACCT CCACCCCCAC CAGACGATGA TCCCACACAC

401 GCTCACCCTG GGGCCCCCCT CCCAGGTCGT CATGCAATAC GCCGACTCCG
 / Pre-1

451 GCAGCCACTT TGTCCCTCGG GAGGCCACCA AGAAAGCCGA GAGCAGCCGG

501 CTGCAG

Fig. 3

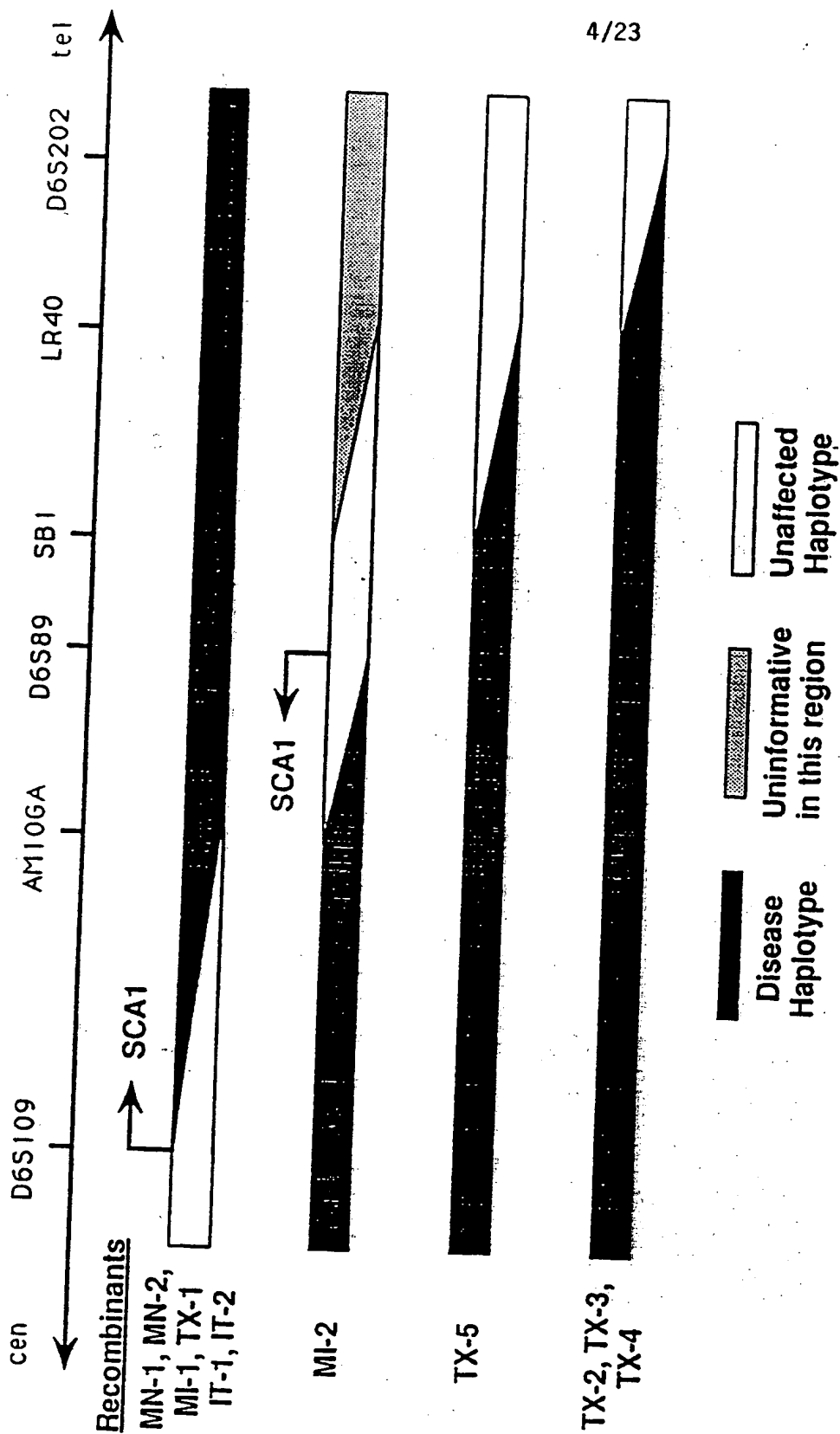


Fig. 4

5/23

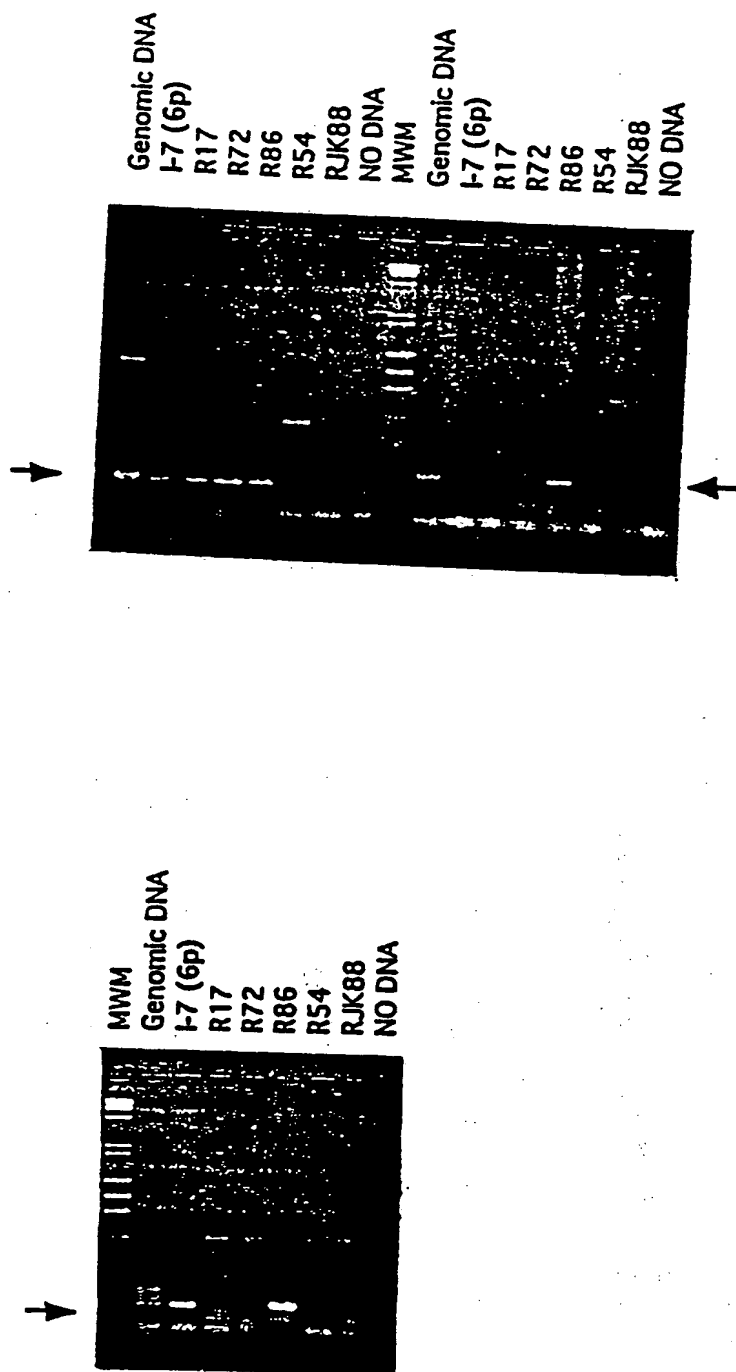


Fig. 5

6/23

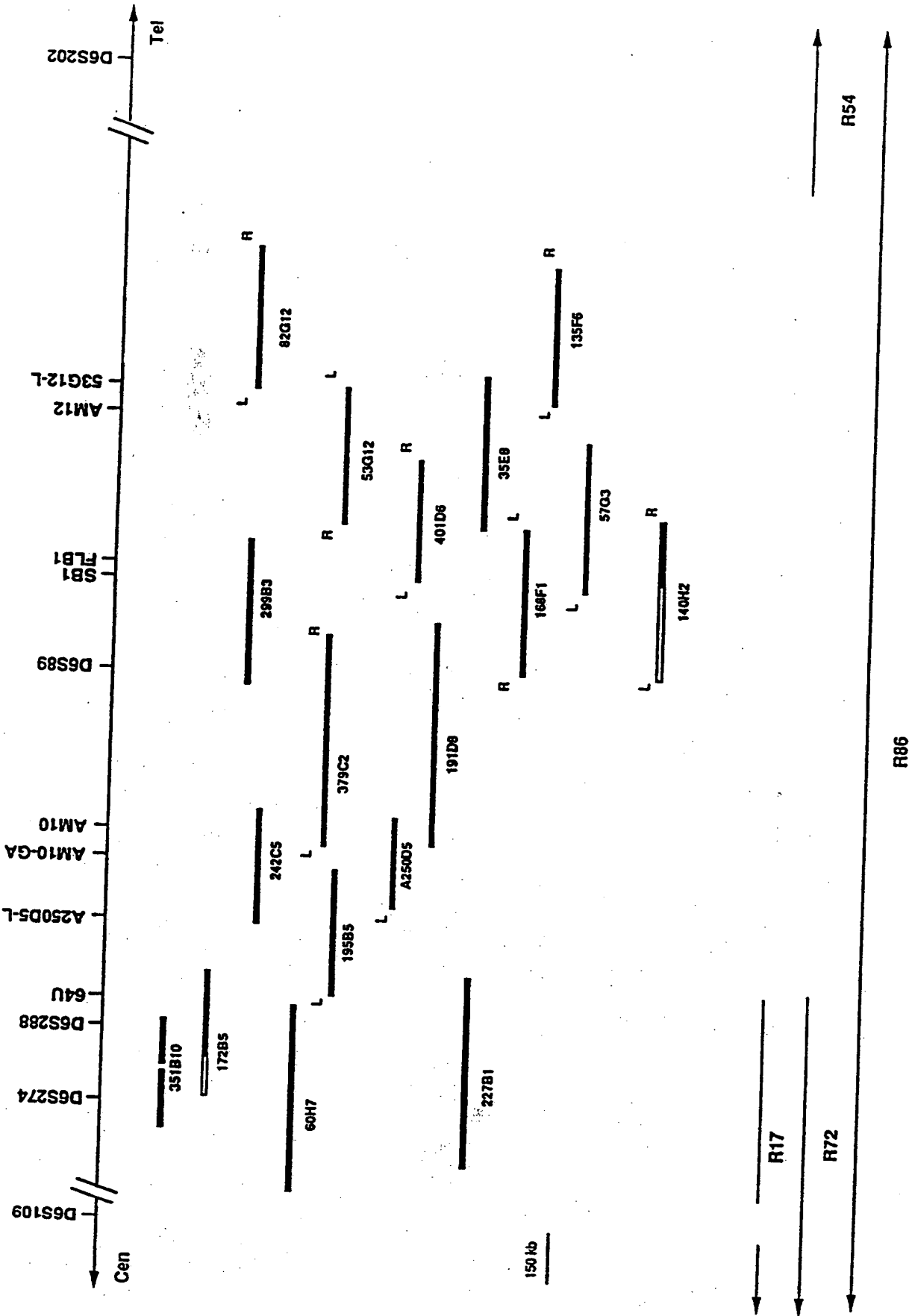


Fig. 6

7/23

SCA1/D6S274 RECOMBINATION EVENT

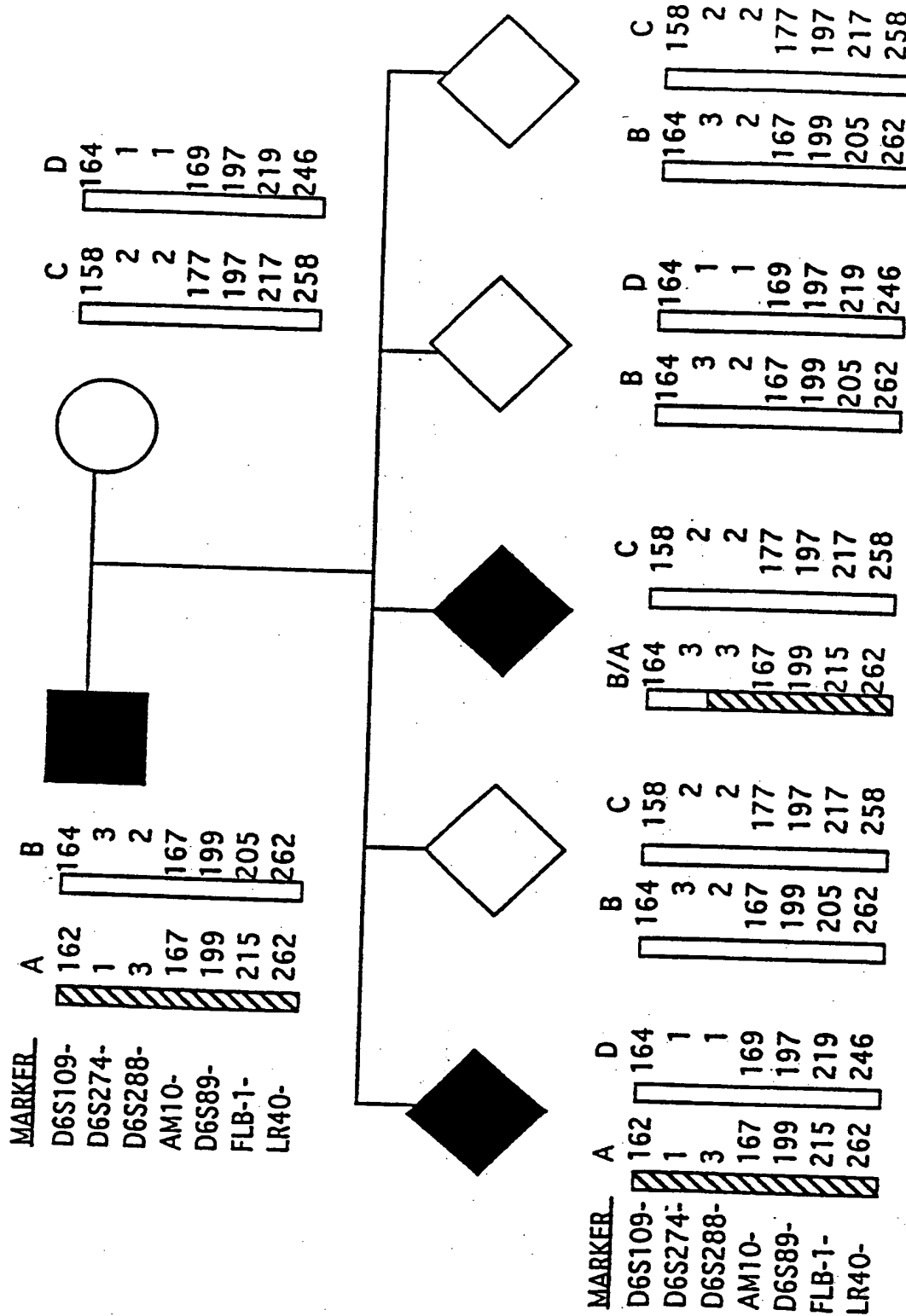


Figure 7

8/23

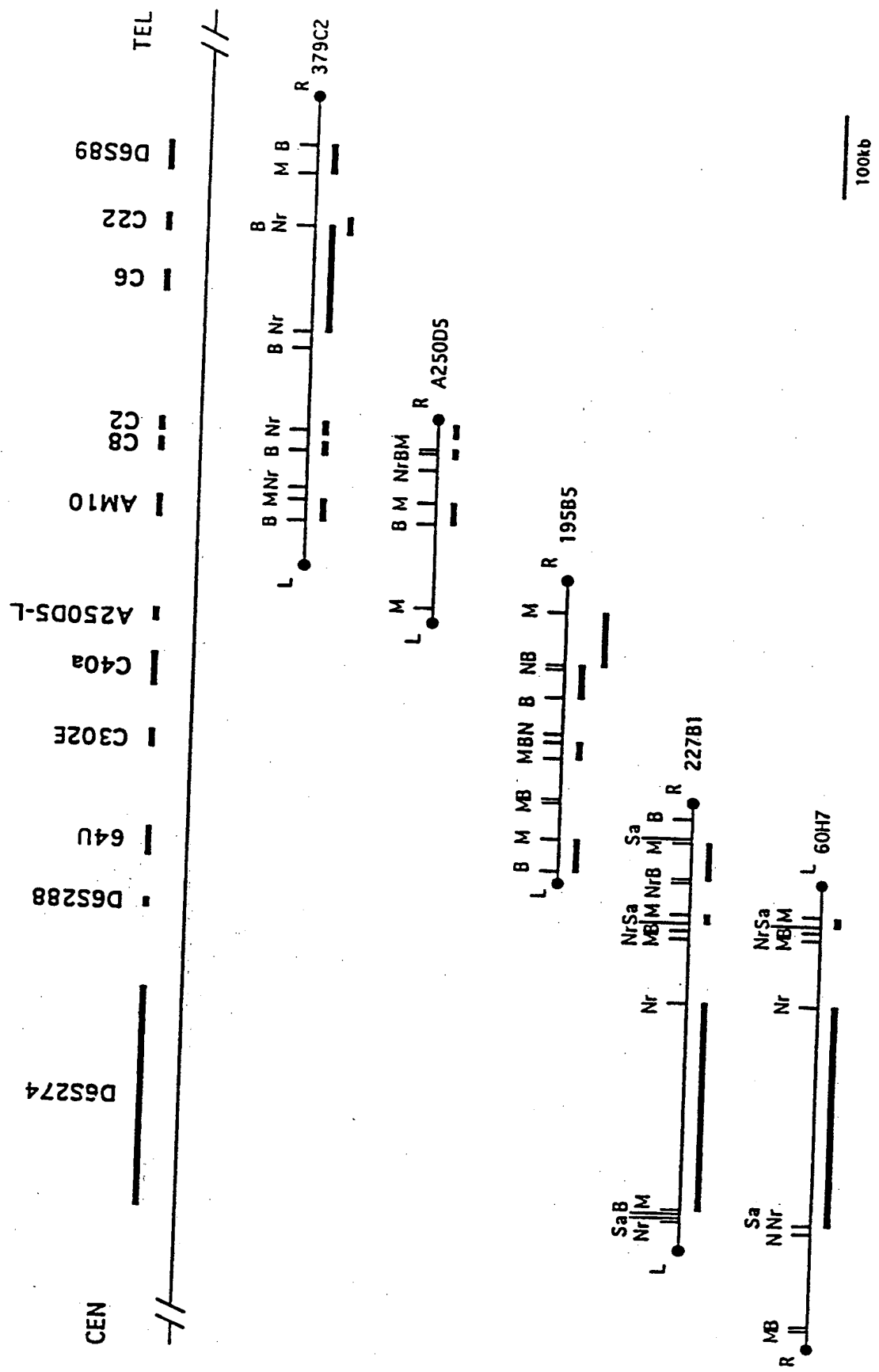


Fig. 8

9/23

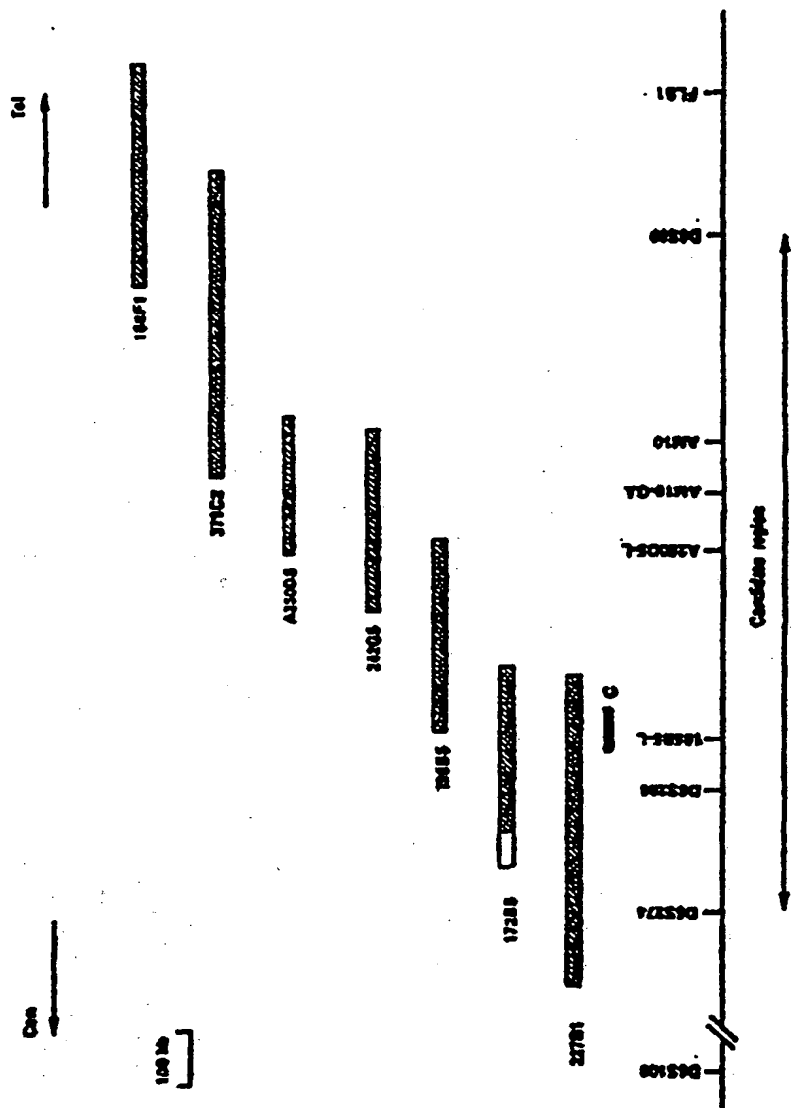


Fig. 9

10/23

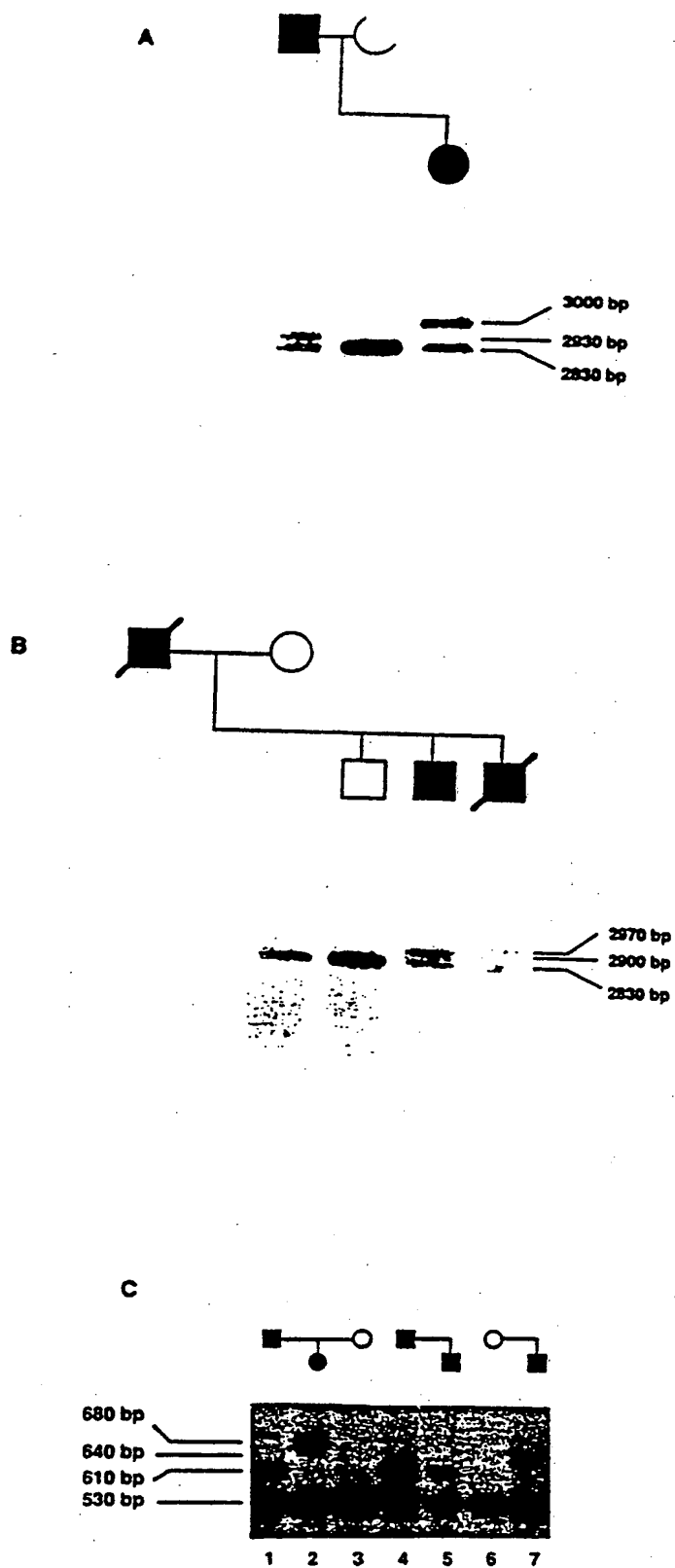


Fig. 10

11/23

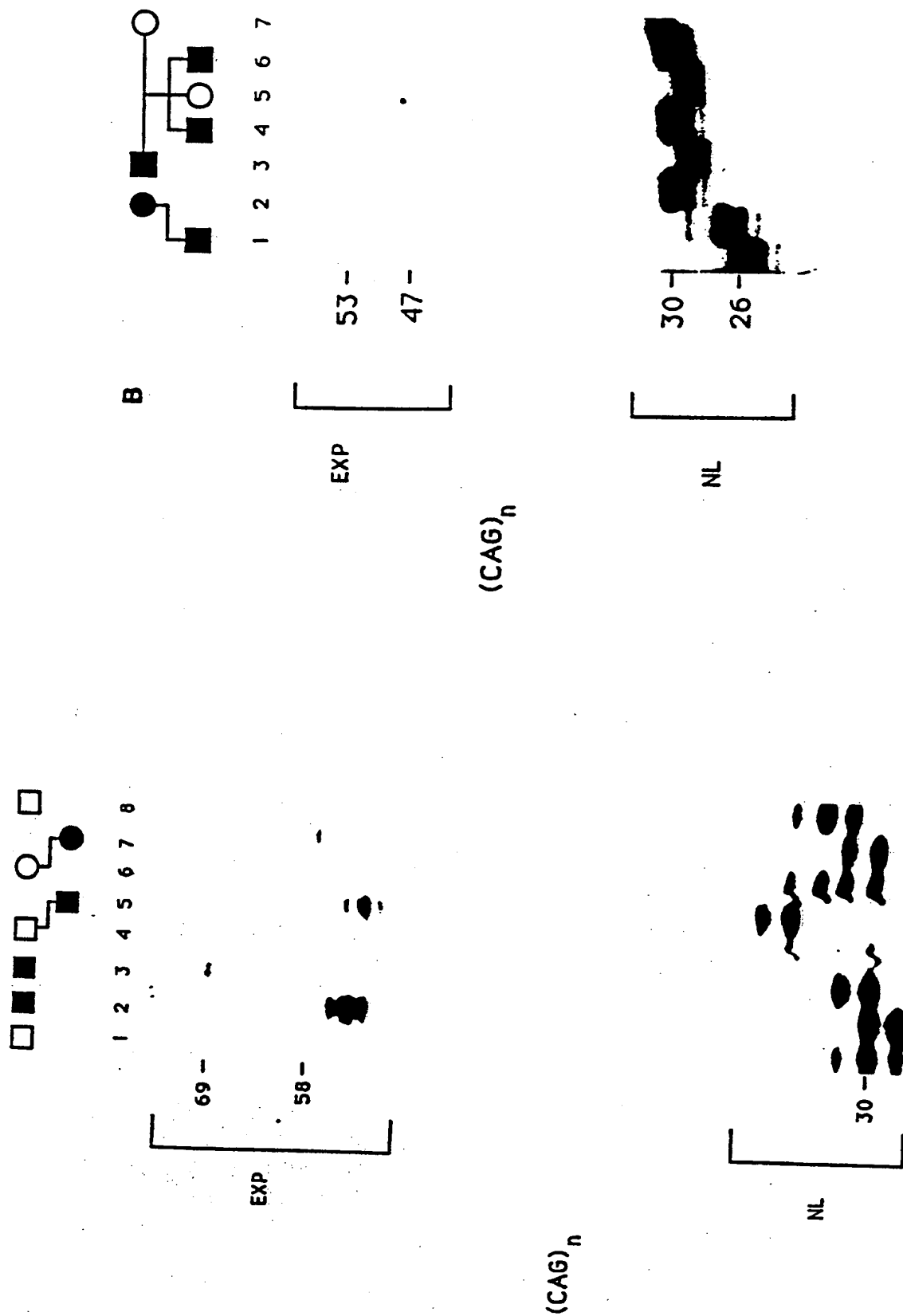


Fig. 11

12/23

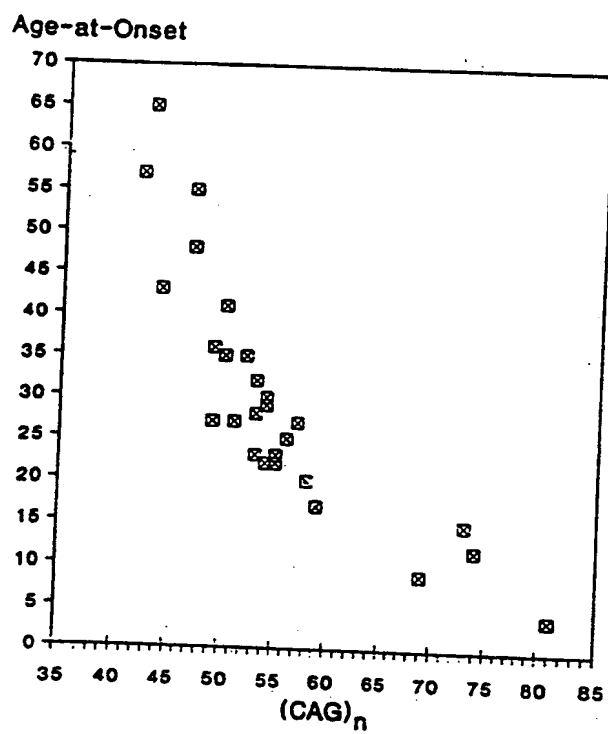


Fig. 12

13/23

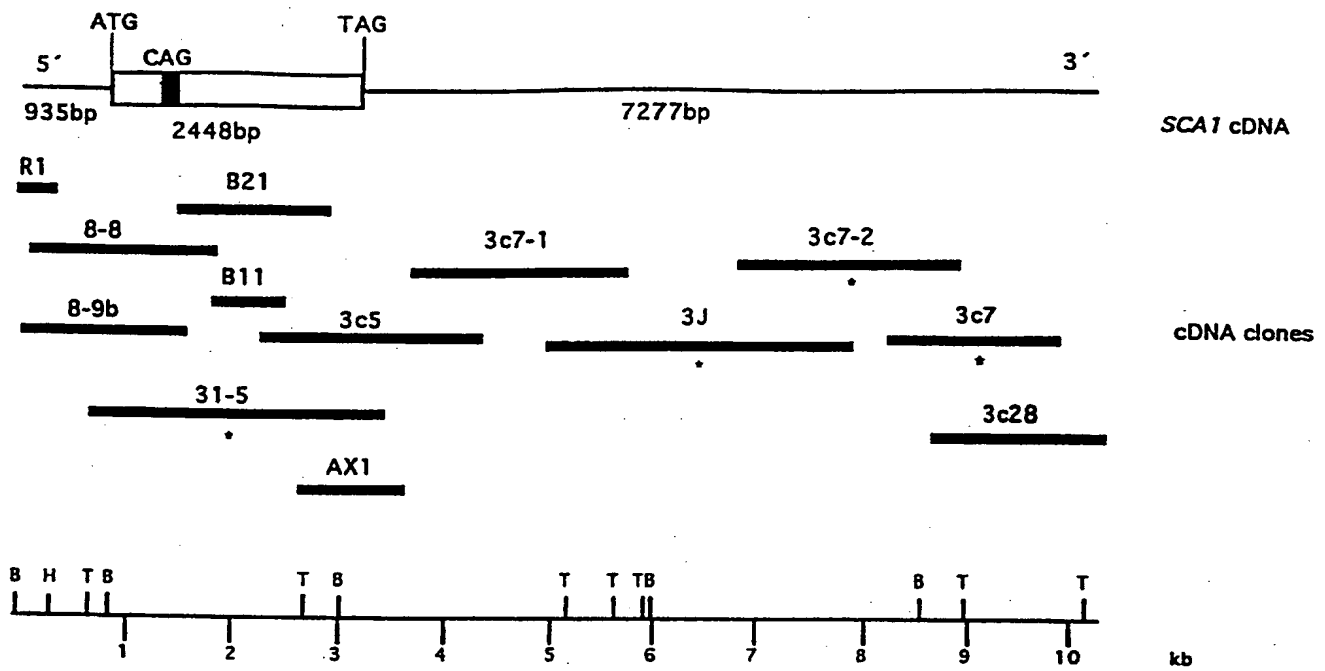


FIG. 13

14/23

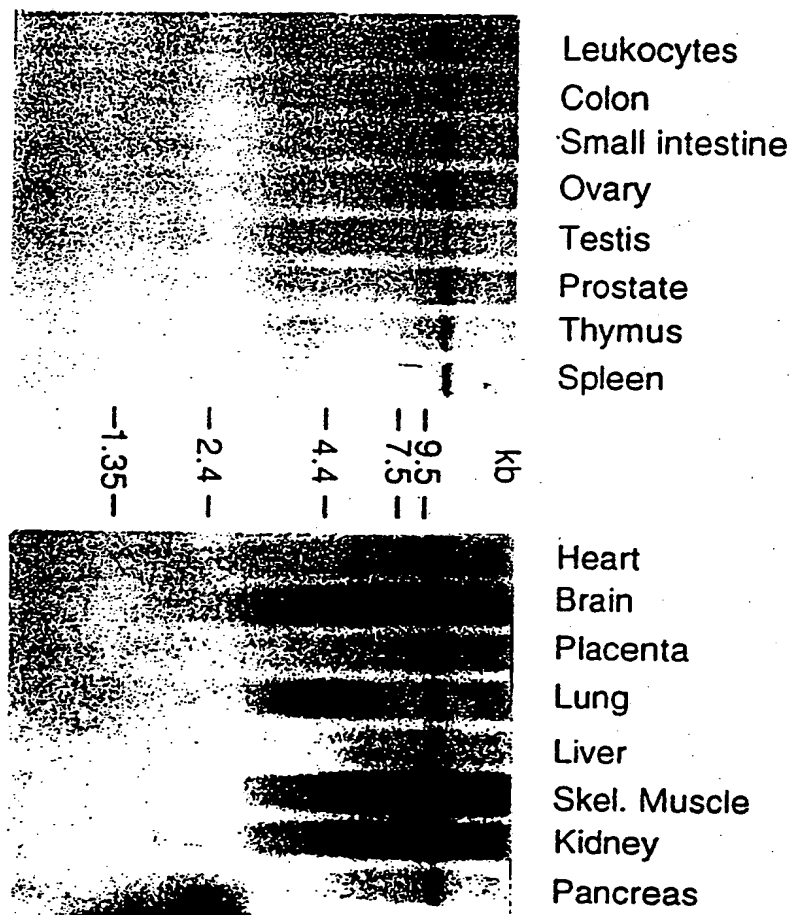


FIG. 14

15/23

FIGURE 15

1 CTACTACAGTGGCGGACGTACAGGACCTGTTTCACTGCAGGGGGATCCAAAACAAGCCCCGTGGAGCAACAGCCAGAGCAACAGCAGCTG 90
91 CAAGACATTGTTTCTCTCCCTCTGCCCCCTTCCCCACGGCAACCCAGATCCATTACACTTTACAGTTTACCTCACAACAACTACTA 180
181 CAAAGCACCAGCTCCCTGATGGAAAGGAGCATCGTGATCAAGTCACCAAGGGTGGTCCATTCAAGCTGCAGATTGTTTGTTCATCCTTGT 270
271 ACAGCAATCTCCTCCTCCACTGCCACTACAGGGAAGTGCATCATGTGCAGCATACTGGAGCATAGTGAAGAGTCTATTTTGAAGCTTC 360
361 AAACCTTAGTGCTGCTGCAGACCAGGAACAAGAGAGAAAGAGTGGATTTCAGCCTGCACGGATGGTCTTGAAACACAATGGTTTTTGGTC 450
451 TAGGCGTTTTACACTGAGATTCTCCACTGCCACCTTTCTACTCAAGCAAAATCTTCGTGAAAAGATCTGCTGCAGGAACTGATAGCTT 540
541 ATGGTTCTCCATTGTGATGAAGACATGGTACAGTTTCCAAAGAAATTAGACCATTTTCTTCGTGAGAAAGAAATCGACGTGCTGTTT 630
631 TCATAGGGTATTTCTCACTTCTCTGTGAAAGGAAGAAACACGCTGAGCCCAAGAGCCCTCAGAGCCCTCCAGAGCCTGTGGGAAG 720
721 TCTCCATGGTGAAGTATAGGCTGAGGCTACCTGTGAACAGTACGCAGTGAATGTTTCATCCAGAGCTGCTGTTGGCGGATTGTACCCACGG 810
811 GGAGATGATTCTCATGAAGAGCCTGGATCCCTACAGAAATCAATGTGACTTTCGTTTATCAGACTAAAATCAGAGCCATCCAGACA 900
901 GTGAAACAGTCACCGTGGAGGGGGGACGGCGGAAATGAAATCAACCAAGAGCGGAGCAACGAATGCCCTGCCCTCCCAAGAAGCGGAGA 990
1 TCCCCGCCACAGCCGGTCTCCGAGGAGAAGGCCCTACCTGCGCCAGCGACAACCCGGTGGAGGGCACAGCATGGCTCCCCGGCA 1080
991 P A T S R S S E E K A P T L P S D N H R V E G T A W L P G N 49
20 ACCTGCTGGCGGGGCCACGGGGGGGAGGATGGCGGGCAGGACCTCGGTGGAGCTTGGTTTACAACAGGGAATAGGTTTACACA 1170
1081 P G G R G H G G R H G P A G T S V E L G L Q Q G I G L H K 79
50 AAGCATTGTCCACAGCTGGACTACTCCCCCGCAGCCTCCAGGTCTGTCCCGTGGCCACCAAGCTGCTGCTGCCGCTAGCCACCC 1260
1171 A L S T G L D Y S P P S A P R S V P V A T T L P A A Y A T P 109
80 CGCAGCCAGGACCCCGGTGTCCTCCCGTGCAGTACGCTCACCTGCCGACACCTTCCAGTTCATTGGGTCTCTCCAATACAGTGAACCT 1350
1261 Q P G T P V S P V Q Y A H L P H T F Q F I G S S Q Y S G T Y 139
110 ATGCCAGCTTCATCCCATCAGCTGATCCCCCAACCGCAACCCCGTCCAGTGCAGTGGCCTCGGCCGAGGGGCCACCATCCAT 1440
1351 A S F I P S Q L I P P T A N P V T S A V A S A A G A T T P S 169
140 CCCAGCGCTCCAGCTGGAGGCTTATCCACTCTGCTGGCAACATGGCAGTCTGAGCCAGACGCGGGACACAAGGCTGAGCAGCAGC 1530
1441 Q R S Q L E A Y S T L L A N M G S L S Q T P G H K A E Q Q Q 199
170 AGC 1620
1531 Q Q Q Q Q Q Q Q Q H Q H Q 229
200 GCAGGGCTCCGGGGCTCATACCCCGGGTCCCCCCCCACCGCCAGCAGAGCAACCAAGTACGTCCACATTCCAGTTCTCCGAGAACCCG 1710
1621 R A P G L I T P G S P P P A Q Q N Q Y V H I S S S P Q N T G 259
230 GCCGCACCGCCTCTCTCCGGCTTCCCGTCCACCTCCACCCCCACCGAGCATGATCCCAACACGCTCACCTGGGGCCCCCTCC 1800
1711 R T A S P P A I P V H L H P H Q T M I P H T L T L G P P S Q 289
260 AGGTCGTATGCAATACGCGGACTCCGGCAGCCACTTGTCTCCCTCGGGAGGCCACCAAGAAAGCTGAGAGCAGCCGGCTGCAGCAGGCCA 1890
1801 V V M Q A D S G S H F V P R E A T T K K A E S S R L Q Q A I 319
290 TCCAGGCCAAGGAGTCTGAACGGTGAAGTGAAGAGAGCCGGTACGGTACGGTCCCTCCTCAGCCGACCTGGCGCTGGGAGGAG 1980
1891 Q A K E V L N G E M E K S R R Y G A P S S A D L G L G K A G 349
320 GCGCAAGTCGGTCTCTCACCCTACGAGTCCAGGACGCTGGTGTCCACCCGAGCCCTCAGACTACAGCAGTCTGATCTCTCGGGG 2070
1981 G K S V P H P Y E S R H V V G H P S P S D Y S S R D P S G V 379
350 TCCGGGCTCTGTGATGCTCTGCCAACAGCAACAGCCCGCAGCTGACCTGGAGGTGCAACAGGCCACTCATCGTGAAGGCTCCGCTT 2160
2071 R A S V M V L P N S N T P A A D L E V Q Q A T H R E A S P S 409
380 CTACCTCAACGACAAAGTGGCTGCAATTAGGGAAGCCTGGCCACCGTCTACCGCTCTCACCCACACGGTCACTCAGACACAC 2250
2161 T L N D K S G L H L G K P G H R S Y A L S P H T V I Q T H 439
410 ACAGTCTCAGAGCACTCCCGTGGGACTGCCAGCCACGGCTTCTACGAGGAACTCAACCCCTGTCTATCGGCTACCTAGCGGCC 2340
2251 S A S E P L P V G L P A T A F Y A G T Q P P V I G Y L S G Q 469
440 AGCAGCAAGCAATCACTACGCGGCTGCCCGGACCTGCTGATCCCGGACACAGCCCTGCTCATCCCGTGGGAGCAGT 2430
2341 Q Q A I T Y A G S L P Q H L V I P G T Q P L L I P V G S T D 499
470 ACATGGAAGCGTCCGGGGCAGCCCGGCTAGTACGCTCATCCCCAGTTTGTGTCAGTGCTCAGCGTTGTCACACGTTGTCACACCGCCCTT 2520
2431 M E A S A P A I V T S S P Q F A A V P H T F V T T A L P 529
500 CCAAGAGCGAGAATCTCAACCTGAGGCCCTGTGTACCCAGGCGCTACCCAGGCTGGTGCAGGCCAGTCCACCTGCTGTGGTGC 2610
529 K S E N F N P E A L V T Q A A Y P A M V Q A Q I H L P V V Q 559
2611 AGTCCGTGGCTCCCCGGGGCGGCTCCCCCTACGCTGCTCCTTCTCATGAAAGGCTCCATCATCCAGTTGGCCAACGGGAGCTAA 2700
560 S V A S P A A A P P T L P P V H P K G S I I Q L A N G E L K 589
2701 AGAAGGTGGAAGACTTAAACAGAAAGATTTCATCCAGAGTGCAGAGATAAGCAACGACCTGAAGATCGACTCCAGCACCGTAGAGAGGA 2790
590 K V E D L K T E D F I Q S A E I S N D L K I D S S T V E R I 619
2791 TTGAAGACAGCCATAGCCGGCGTGGCCGTGATACAGTTCGCGTGGGGAGCAGCCAGCCAGGTGAGGTTTGGTAGAGT 2880
620 E D S H S P G V A V I Q F A V G E H R A Q V S V E V E Y 649
2881 ATCCTTTTTTGTGTTTGGACAGGGCTGGTCTCTGCTGTCCGAGAGAACAGCCAGCTCTTTGATTGCGGTGTTTCAAACCTCTCAG 2970
650 P F F V F G Q G W S S C C P E R T S S Q L F D L P C S K L S V 679
2971 TTGGGATGCTGCTGCTCGCTTACCCTCAAGAACTGAAGAACGGCTCTGTAAAAAGGGCCAGCCCGTGGATCCCGCCAGCGTCTG 3060
680 G D V C I S L T L K N L K N G S V K K G Q P V D P A S V L L 709
3061 TGAAGCACTCAAAGGCCAGGCGCTGGCGGGCAGCAGACAGGTATGCCGAGCAGGAAACGGAATCAACAGGGGAGTGCCAGATGC 3150
710 K H S K A D G L A G S R H R Y A E Q E N G I N Q G S A Q M L 739
3151 TCTCTGAGAATGGCGAAGTGAAGTTTCCAGAGAAATGGGATTGCTGAGCGCCCTTCTCACCAGAAATAGAACCCAGCAGCCCGCG 3240
740 S E N G E L K F P E K M G L P A A P F L T K I E P S K P A A 769
3241 CAACGAGGAAGAGGAGGTGGTGGCGCCAGAGACCGCAACCTGGAGAGTGCAGAGACGAACCACTTGAAGTCTTCTAAGCTTCTC 3330
770 T R K R W S A P E S R K L E K S E D E P P L T L P K P S L 799
3331 TAATTCCTCAGGAGGTTAAGATTGTCATTGAAGGCCGCTCTAATGTAGGCAAGTAGAGGAGCGTGGGGGAAAGGAACTGGCTCTCCC 3420
800 I P Q E V K I C I E G R S N V G K * 829
3421 TTATCATTGTATCCAGATTACTGTACTGTAGGCTAAAATAACACAGTATTTACATGTTATCTTCTAATTTTAGGTTTCTGTTCTAACC 3510
3511 TTGTCATTAGATTACAGCAGGTGTGTCGAGGAGACTGGTGCATATGCTTTTCCACAGTGTCTGTGAGTGGCGGGGGAGGAAGG 3600

3690	GCACAGCAGGAGCGGT CAGGGCTCT CACGGCATCCCCGGGGGAGAAAGGAACGGGGCTT CACAGTGCCTGCCTTCTCTAGCGGCACAGAAGC	3690
3691	AGCCGGGGGCGCTGACTCCCGCTAGTGT CAGGAGAAAAAGTCCCGTGGGAAGAGTCTCTGCGGGGTG CAGGGGTGCACGCATGTGGGGGTG	3780
3781	CACAGGCGCTGTGGCGCGAGTGGAGGCTCTCTTTTTCTGCGTCCCTCTGCTCTCTCTGCTATCGGCTAGGCTGGCGGGGGGTTCA	3870
3871	GAGCAGTGTCTCTCTGGGGTCCCGAGTGCACAAATCAACATCAGGAACCCAGCTTCAGGGCATCGCGGAGACGCGT CAGATGGCAGATTT	3960
3961	GGAAAGTTAACCATTTAAAAGAACATTTTTCTCTCCAACATATTTTACAATAAAAGCAACTTTTAATTGTATAGATATATATTTCCCCCT	4050
4051	ATGGGGCTGACTGCACGTATATATATTTTTTTAAAGAGCAACTGCCACATGCGGGATTTCAATTTCTGCTTTTTTACTAGTGCAGCGATG	4140
4141	TCACCAGGGTGTGTGGTGGACAGGGAAGCCCTGCTGTCTAGTGGCCACATGGGGTAAGGGGGTGTGGGGTGGGGGAGGAGGAGAG	4230
4231	CGAACCACCCAGCTGGTTTCTGTGAGTGTGTAGAAACCAATCAGGTTATTGTCATTGACTTCACTCCCAAGAGGTAGATGCAAACTGCC	4320
4321	CTTCAGTGAGAGCAACAGAAGCTCTTACGTTGAGTTTGGGAAATCTTTTGTCTTTGAACTCTAGTACTGTTTATAGTTCATGACTAT	4410
4411	GACAACCTCGGGTCCCACTTTTTTTTTTTCAGATTCCAGTGTGACATGAGGAATTAGATTGGAAGTACGACATATTAATCTATCTTTAA	4500
4501	GCATTTAAAAATACGTTGCACACTTTATACCAAGCATCTTGGTCTCTCATCAACAGTACTGATCTCACTTTAAACTCTTTGGGGAA	4590
4591	AAAACAAAAACAAAAAACTAAGTTGCTTTCTTTTTTCAACACTGTAACATCAATTCAGCTCTGCAGAAATGCTGAAGAGCAAGATAT	4680
4681	TGAAAGTTTCAATGTGGTTTAAAGGGATGAATGTGAATTATGAACCTAGTATGTGACAAATAAATGACCACCAAGTACTACCTGACGGGAG	4770
4771	CACTTTTCTACTTTGATGTCTGAGAACTCAGTTAGAGGCATATG CAGAGTTGGCAGAGAACTGAGAGAAAGAGGATGGAGAGAAATACT	4860
4861	CATTTTTGTCTGAGTGTTTTTCTTTTAAAGTGAACTTTTAAAGAACCTTGGCATTTGCACATATTGAGTTTATAACTTGTGTGATATTC	4950
4951	TGCAGTTTTTATCCAATAACATTGTGGGAAAGGTTTGGGGGACTGAACGAGCATAAATAAATGTAGCAAAATTTCTTTCTAACCTGCCTA	5040
5041	AACCTAGGCACTTTTATAAGGTTATGTTCTTTGAAATCATTTTGGTCTTTTACCACATCTGTCAACAAAAAGCCAGGTTCTTAGCGG	5130
5131	GCTCTTAGAAACTCTGAGAATTTTCTTCAGATTCAATGAGAGAGTTTCCATAAAGACATTTATATGTGAGCAGATTTTTTTTAAAC	5220
5221	AATTAATCTTATATGTGTGTTAATAGTTATTTTTCAGAAAGGCTTTTTTTTCTATTCAAATCAAATCGAGATTTAATGTTTGGTACA	5310
5311	AACCCAGAAAGGGTATTTCATAGTTTTTAAACCTTTCATTTCCAGAGACTCCGAAATATCATTTGTGGGTTTTGAATGCATCTTTAAAGT	5400
5401	CTTTAAAAAAAGGTTTTTATAAGTAGGAGCAAAATTTTTTAAATATTTCTTGTGATGGCTGCACTAACTGAACAAATACCTGACTTTTC	5490
5491	TTTTACCCCATTTGAAATAGTACTTTCTCGTTTTCAAAATAAAAAAAACTGGTATCAACCCACATTTTGGCTGTCTAGTATTCAT	5580
5581	TTACATTTAGGGTTTACCAGGACTAATGATTTTTATAAACCGTTTTCTGGGGTGACCAAAAAATTGAAATAGGTTTAGAATAGCTAGA	5670
5671	ATAGTCTCTGACTTTCTCGAATTTCAATACCCTCTCAGACATGCTTGCAGAGAGCTGGGGGGCTCTTGTGACTTCTGCACTACTGCTTAT	5760
5761	TTAGTGTCTGATTTTTTAAACGTTTCTGTT CAGAGAACTTGCTTAATCTTCCATATATCTGCTCAGGGCTGCAATTAATAGGTTTT	5850
5851	GTTTTTCTTTTGTGTTTTTGTGCTTTGAGCTTTGATGGTAAAGAGGAATACGGGCTGCCCATAGACTTTGTTCTCATTAAATCACTATTTCAACT	5940
5941	CATGTGGACTCAGAAAAACACACACCCTTTTGGCTTACTTCGAGTATGAATTGACTGGATCCACTAAACCAACTAAGATGGGAAA	6030
6031	ACACACATGTTTGGAGCAATAGGAACATCATCAATTTTTGTGGTCTTATTTCCAGTATAGGAATTAAAAATTAGTTCTTTCTTA	6120
6121	AACACTTGTCCATTTCAATCTCTGCTTTTTTAGCATGTGCAATTTCTTCTGCGCAATAGAGTCTGACCACTGTGCTATATTTTAA	6210
6211	GCTCATTCCCTTTTGGCTTTTTCTGTTTGGTTGATCTTCCCATTTCTGGCCAGAGCAGGGCTGAGGGAAGGAGCCAGGAGGGAGAGA	6300
6301	GCCTCCACCTTTCCCTGCTCGGGAGTCTGAGTGTGGGGCGGGGAGCCCTCAGGAGGCCCTGCGCTGCGCCAGCTTGCAGGAAGA	6390
6391	GCCAGCCAAAGGAGACCCGGGGGAGGAACCCGCAAGTGTCCCTGT CACCACGGAATAGTGAATGTGGAGTGTGGAGAGGAAGGAGCAGA	6480
6481	TTCAATTTAAGACGCACTCTGGAGCCATGTAGCTCTGGAGTCAACCCATTTTCCACGGTCTTTTCTGCAAGTGGGCAGGCCCCCTCTCG	6570
6571	GGTCTGTGCTCTTGAGACTTGGAGCCCTGCTCTGAGCCTGGACGGGAAGTGTGGCCGTGTGTGTGTGCGTTCTGAGCGTCTTGGCCA	6660
6661	GTGGCTGTGGAGGGGACCACTGCGCCACCCACCGTCAACCTCCCTTCTGTCGGCAGCTTTCTCTCAAATAGGAAGAACGCACAGAGGCGAG	6750
6751	AGCCTCTGTTTTCAGACGTTGGCGGCCCGGAGGCTCCAGAGAGCTGCTGTCAACCGCTCTGTGTGTAGCAAACTTAACGATGACAGG	6840
6841	GGTAGAAATTTCTCGGTGCGGTT CAGCTTACAAGGATCAGCCATGTGCTCTGTACTATGTCCAATTTGCAATATTACCGACAGCCGTC	6930
6931	TTTTGTTCTTTCTTTCTGTTTTTCCATTTTTTAAACTAGTAAACAGCGGCCCTTTGCGTTTTACAAATGGAACCAATACCAAGAAATTAGT	7020
7021	CAGGGCGAAAGAAAAATAATCATTTAATAAGAAACCAACAAACAGAACTCTCTTTCTAGGGATTTCAAATATATAAAATGACT	7110
7111	GTTCTTAGAATGTTTAACTTAAGAATTATTACAGTTTGTCTGGGCCACTGCGGGCAGAGGGGGAGGGAGGGATACAGAGATGGATGC	7200
7201	CACCTACCTCAGATCTTTTAAAGTGGAAATCCAAATGAAATTTTCAATTTGGACTTT CAGGATAAATTTTCTATGTTGGTCAACTTTTCGTT	7290
7291	TTCCCTTAAGTCCACCGGATTTGAGTTTGGGATGATTTGATTTCTGTTGTGTGATGCCATTTCTAATTTGGAATTTGTAGCCTCTATGTT	7380
7381	TCGTTAGGTGAGTGTGTTGGTTTTTTTCCCCACAGGAAGTGGCAGCATCCCTCTCTCCCTTAAAGGAGCTCTGCGAACCTTTTC	7470
7471	ACACCTCTTTCTCAGGAGCGGGGAGGTGTGTGTGGGTACACTGACGTGTCCAGAAGCAGCACTTTGACTGCTCTGGAGTAGGGTTGTA	7560
7561	CAATTTCAAGGAATGTTTGGATTTCTGTCATCTGTGTGATTCTCTTAGATACCGCATAGATTGCAATATATATGTCATGTTCAAGAT	7650
7651	GAAACAGTAGTCTCTAGTAATCATAAATCCACTTTGACAGTTTCTTACAGTAATGTTTCTTACGAAATATGTTGCCAAATTTATTTTGTGTGT	7740
7741	AGCTCTGGAATTTGTTTGTGTTTTTAAAGGAACGATTGACAAATACCTTTTAACTCTGTGACTACTAAGGAACCTATTTCTTTT	7830
7831	ATAGAGAGAAAAATCTCCAATGCTTTTGAAGCACTAATACCGTGTCTTTT CAGATATGGGTGAGGAAGCAGAGCTCTCGGTACCGAAGG	7920
7921	CGGGGCTTTCTGAGCTGTGTGGTTGTCTATGGCTATGTTTCATGAACCAACAGCAGCTCAACAGACTGTGCTGTGCTCTGTTGAAACCC	8010
8011	TTTGCACTTCAATTTGACAGGAGTGAAGAACAGGGCCAGCAGACTCCATGGCCCAATTCGGTTTCTTCTGGTGGTATGTGAAGGAGAGAA	8100
8101	TTACACTTTTTTTTTTTTAAAGTGGCGTGGAGGCCCTTGTCTCCACATTTGTTTTTAAACCGAAATTTCTGAATAGAGAAATTAAGAAC	8190
8191	ACATCAGTAATAATAATACAGAGAAATACTTTTTTATAAAGCACATGCATCTGCTAATGTGTTGGGTGTTTCTCTCTTTTCCAGC	8280
8281	GACAGTGTGTGTTTTCTGGCATAGGGAATCTCAAACACTTGACACCTCTACTCCGAGACTGAGATTTCTTTACATAGATGACACTCG	8370
8371	CTTCAAATACGTTACCTTACTGATGATAGGATCTTTCTGTAGCACTATACCTTGTGGGAATTTTTTTTTTAAATGTACACTGATTGTA	8460
8461	GAGCTGAAGAAAAACAAATTTTGAAGCACTCACTTTGAGGAGTACAGGTAATGTTTAAAAAATTTGCACAAAAAGAAAAATGAATGTGTA	8550
8551	AATGATTCAATCAGTGTTTGAAAGATATGGCTCTGTTGAAACATAGGTTTCACTATTTGTTTGTAA	

[illegible]

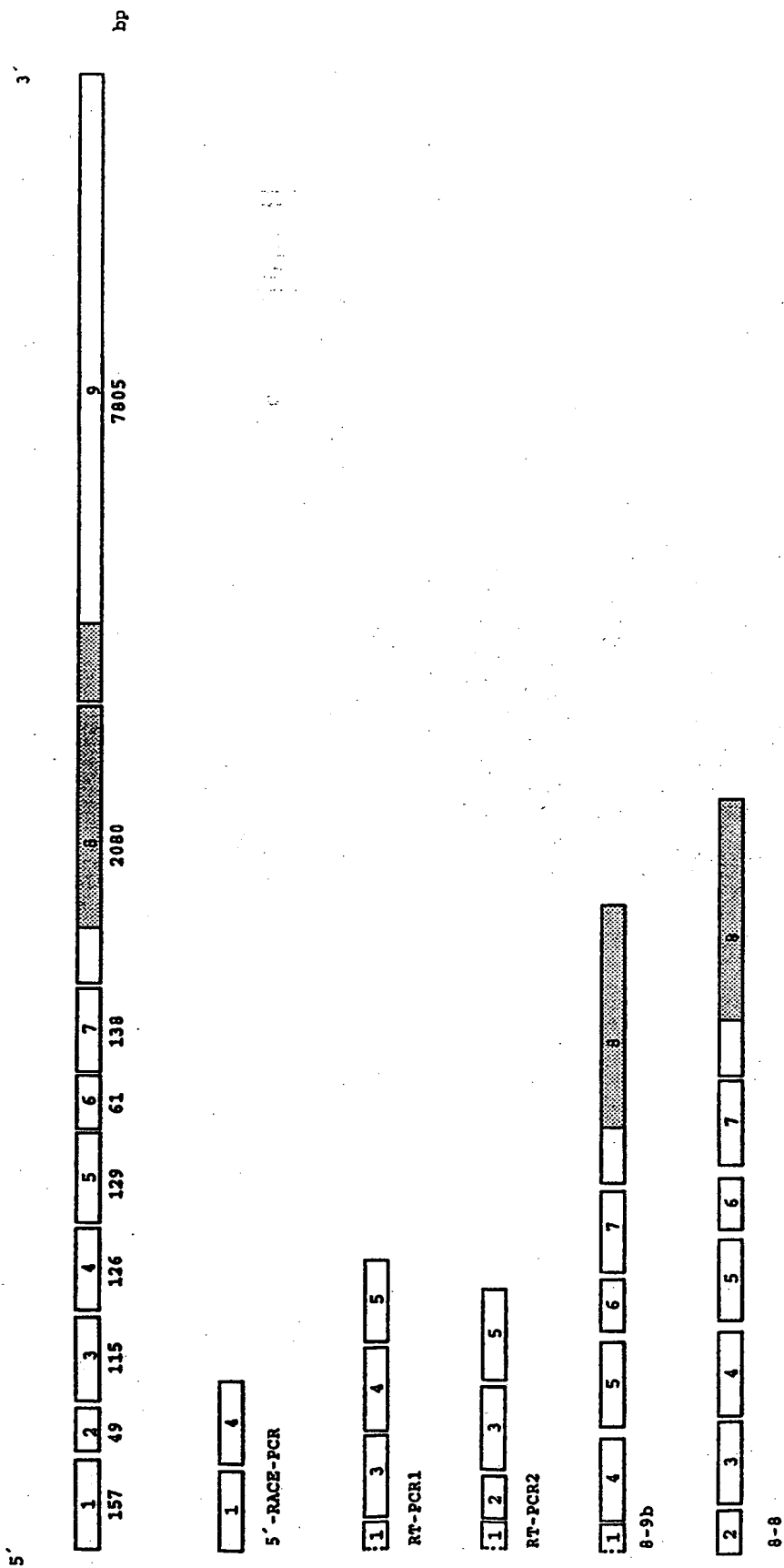


FIG. 16a

19/23

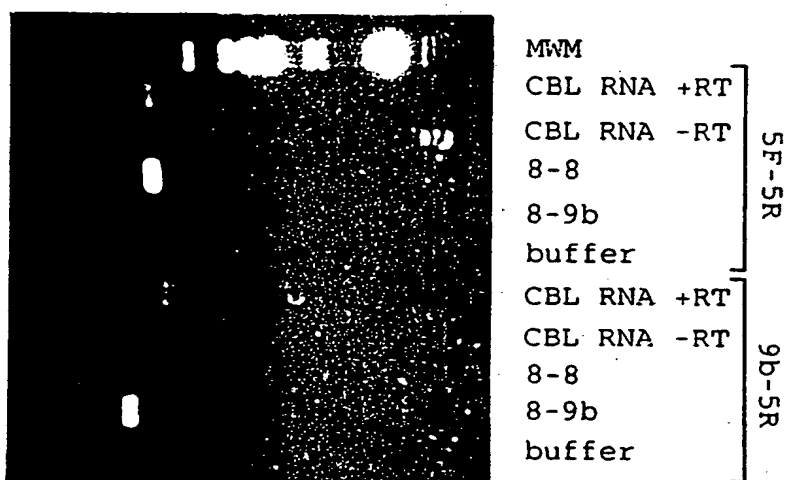


FIG. 16b

20/23

		Exon 1	157 TTTACA gtaagtga
gtttctatgcatag	158 GTTTACC	Exon 2	206 GGAAAG gtatatgg
ctcgaccattgcag	207 GAGCATCG	Exon 3	321 TGTCAG gtgagagt
ttgtttgactgcag	322 CATACTGG	Exon 4	447 TTTTTG gtaagtca
ttttataattacag	448 GTCTAGGC	Exon 5	575 GTACAG gtaaacad
tttttctattccag	576 TTTTCCAA	Exon 6	637 CATAGG gtgagtga
tatttccatgctag	638 GTATTTCT	Exon 7	775 AATGTT gtaagtta
cttccctttccag	776 CATCCAGA	Exon 8	2855 GCCCAG gtaacgtt
ccctgtttccacag	2857 GTCAGCGT	Exon 9	

YYYYYYYYYYNCAG

Consensus

AG GTRAGT

FIG. 17

21/23

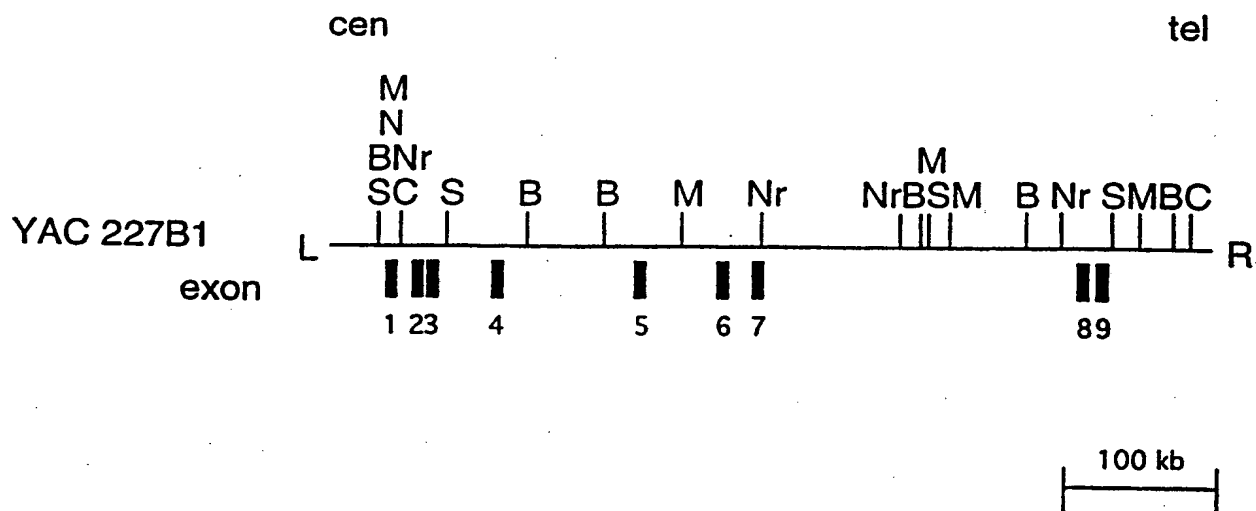


FIG. 18

22/23

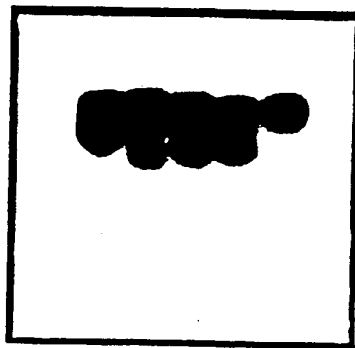


FIG. 19

23/23

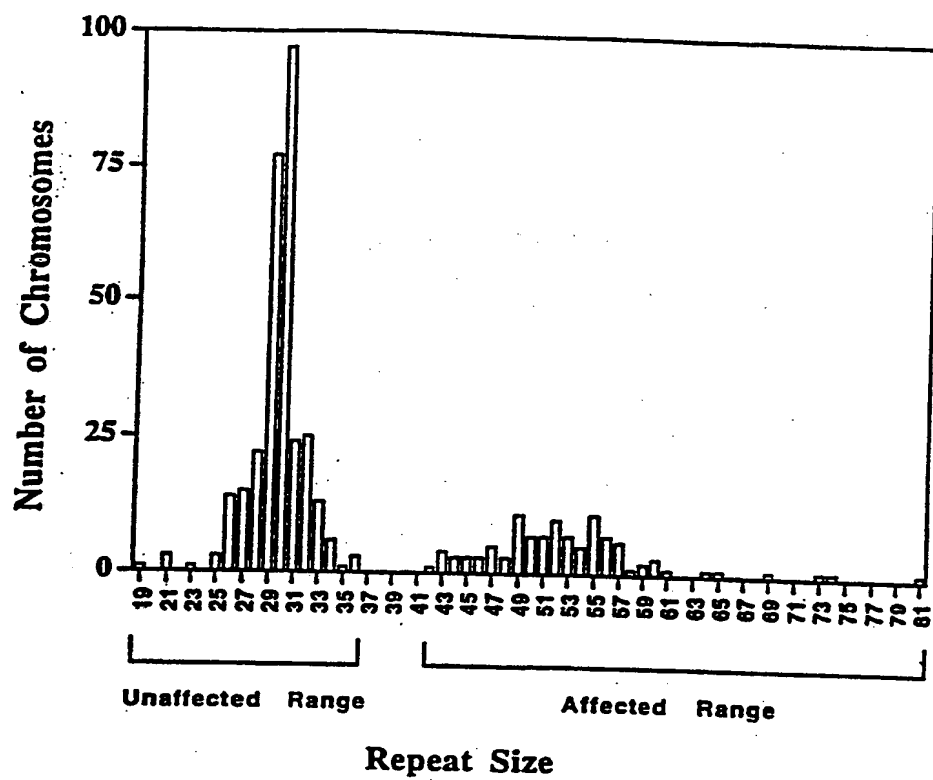


FIG. 20

PCTWORLD INTELLECTUAL PROPERTY ORGANIZATION
International Bureau

INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁶ : A61K 38/00, 15/31, 15/09, 48/00, C12N 15/79, 15/63, 15/00, C07K 16/00, C07H 21/00	A1	(11) International Publication Number: WO 97/18825 (43) International Publication Date: 29 May 1997 (29.05.97)
(21) International Application Number: PCT/US96/18370 (22) International Filing Date: 15 November 1996 (15.11.96) (30) Priority Data: 60/006,882 17 November 1995 (17.11.95) US (71) Applicant (for all designated States except US): THE UNIVERSITY OF BRITISH COLUMBIA [CA/CA]; 2075 Westbrook Mall, Vancouver, British Columbia V6T 1Z1 (CA). (72) Inventors; and (75) Inventors/Applicants (for US only): KALCHMAN, Michael [CA/CA]; #502-2233 Allisson Road, Vancouver, British Columbia V6T 1T7 (CA). HAYDEN, Michael, R. [US/CA]; 4484 West Seventh, Vancouver, British Columbia V6R 1W9 (CA). (74) Agents: LARSON, Marina et al.; Oppedahl & Larson, Suite 309, 1992 Commerce Street, Yorktown Heights, NY 10598-4412 (US).		(81) Designated States: CA, JP, US, European patent (AT, BE, CH, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE). Published <i>With international search report.</i>
(54) Title: PROTEIN WHICH INTERACTS WITH THE HUNTINGTON'S DISEASE GENE PRODUCT, cDNA CODING THEREFOR, AND ANTIBODIES THERETO (57) Abstract A protein, designated as HIP1, interacts differently with the gene product of a normal (16 CAG repeat) and an expanded (> 44 CAG repeat) HD gene. The HIP1 protein originally isolated from the yeast two-hybrid screen is encoded by a 1.2 kb cDNA, devoid of stop codons, that is expressed as a 400 amino acid polypeptide. By further screening of a human frontal cortex cDNA library, and employing the protocol for 5' Rapid Amplification of cDNA ends (RACE), a total of 4795 nucleotides (with an open reading frame of 914 amino acids) of the 10 kb message HIP1 have been isolated to date. Expression of the HIP1 protein was found to be limited to the brain, where the interaction of the HIP1 with the HD protein appears to be necessary for the association of the HD protein with the membrane or specific cytoskeletal components to render it functional. Because HIP1 interacts with expanded HD protein less well than with normal length HD, introduction of additional HIP1 or overexpression of HIP-1 can lead to increased functionality of the defective or normal HD protein. Alternatively, modified forms of the HIP1 which bind more effectively to expanded HD could be introduced to convert the expanded HD protein into a functional molecule.		

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AM	Armenia	GB	United Kingdom	MW	Malawi
AT	Austria	GE	Georgia	MX	Mexico
AU	Australia	GN	Guinea	NE	Niger
BB	Barbados	GR	Greece	NL	Netherlands
BE	Belgium	HU	Hungary	NO	Norway
BF	Burkina Faso	IE	Ireland	NZ	New Zealand
BG	Bulgaria	IT	Italy	PL	Poland
BJ	Benin	JP	Japan	PT	Portugal
BR	Brazil	KE	Kenya	RO	Romania
BY	Belarus	KG	Kyrgyzstan	RU	Russian Federation
CA	Canada	KP	Democratic People's Republic of Korea	SD	Sudan
CF	Central African Republic	KR	Republic of Korea	SE	Sweden
CG	Congo	KZ	Kazakhstan	SG	Singapore
CH	Switzerland	LI	Liechtenstein	SI	Slovenia
CI	Côte d'Ivoire	LK	Sri Lanka	SK	Slovakia
CM	Cameroon	LR	Liberia	SN	Senegal
CN	China	LT	Lithuania	SZ	Swaziland
CS	Czechoslovakia	LU	Luxembourg	TD	Chad
CZ	Czech Republic	LV	Latvia	TG	Togo
DE	Germany	MC	Monaco	TJ	Tajikistan
DK	Denmark	MD	Republic of Moldova	TT	Trinidad and Tobago
EE	Estonia	MG	Madagascar	UA	Ukraine
ES	Spain	ML	Mali	UG	Uganda
FI	Finland	MN	Mongolia	US	United States of America
FR	France	MR	Mauritania	UZ	Uzbekistan
GA	Gabon			VN	Viet Nam

- 1 -

PROTEIN WHICH INTERACTS WITH THE HUNTINGTON'S DISEASE GENE
PRODUCT, cDNA CODING THEREFOR, AND ANTIBODIES THERETO

BACKGROUND OF THE INVENTION

This application relates to a protein designated as HIP1 which interacts with the Huntington's Disease gene product, cDNA coding for HIP1, and methods and compositions relating thereto.

5 "Interacting proteins" are proteins which associate *in vivo* to form specific stable complexes. Non-covalent bonds, including hydrogen bonds, hydrophobic interactions and other molecular associations form between the proteins when two protein surfaces are matched or have affinity for each other. This affinity or match is required for the recognition of the two proteins, and the formation of a stable interaction. Protein-protein interactions are
10 involved in the assembly of enzyme subunits; in antigen-antibody reactions; in forming the supramolecular structures of ribosomes, filaments, and viruses; in transport; and in the interaction of receptors on a cell with growth factors and hormones.

Huntington's disease is an adult onset disorder characterized by selective neuronal loss in discrete regions of the brain and spinal chord that lead to progressive
15 movement disorder, personality change and intellectual decline. From onset, which generally occurs around age 40, the disease progresses with worsening symptoms, ending in death approximately 18 years after onset.

The biochemical cause of Huntington's disease has thus far not been determined. Various theories have been advanced, but each has failed to stand up to
20 experimental evidence designed to test its validity. For example, it was suggested that the selective neuronal loss could be attributed to restricted expression of mRNA or proteins in cells undergoing degeneration. No obviously altered levels of mRNA transcript or protein expression has ever been observed in HD-affected tissues, however.

While the biochemical cause of Huntington's disease has remained elusive, a
25 mutation in a gene within chromosome 4p16.3 subband has been identified and linked to the disease. This gene, referred to as the Huntington's Disease or HD gene, contains three repeat regions, a CAG repeat region and two CCG repeat regions. Testing of Huntington's disease patients has shown that the CAG region is highly polymorphic, and that the number of CAG

- 2 -

repeat units in the CAG repeat region is a very reliable diagnostic indicator of having inherited the gene for Huntington's disease. Thus, in control individuals and in individuals suffering from neuropsychiatric disorders other than Huntington's disease, the number of CAG repeats is between 9 and 35, while in individuals suffering from Huntington's disease the number of CAG repeats is expanded and is 36 or greater.

The protein product encoded by the HD gene has been localized to the cytoplasm, including to the membranes of vesicles on the brain of both normal and HD-affected individuals. To date, no differences have been observed at either the total RNA, mRNA or protein levels between normal and HD-affected individuals. Thus, the function of the HD protein and its role in the pathogenesis of Huntington's Disease remain to be elucidated.

SUMMARY OF THE INVENTION

We have now identified a protein, designated as HIP1, that interacts differently with the gene product of a normal (16 CAG repeat) and an expanded (>44 CAG repeat) HD gene. The HIP1 protein originally isolated from the yeast two-hybrid screen is encoded by a 1.2 kb cDNA, devoid of stop codons, that is expressed as a 400 amino acid polypeptide. By further screening of a human frontal cortex cDNA library, and employing the protocol for 5' Rapid Amplification of cDNA ends (RACE), a total of 4795 nucleotides (with an open reading frame of 914 amino acids) of the 10 kb message HIP1 have been isolated to date. Expression of the HIP1 protein was found to be limited to the brain, where the interaction of the HIP1 with the HD protein appears to be necessary for the association of the HD protein with the membrane or specific cytoskeletal components to render it functional. Because HIP1 interacts with expanded HD protein less well than with normal length HD, introduction of additional HIP1 or overexpression of HIP-1 can lead to increased functionality of the defective or normal HD protein. Alternatively, modified forms of the HIP1 which bind more effectively to expanded HD could be introduced to convert the expanded HD protein into a functional molecule.

BRIEF DESCRIPTION OF THE DRAWING

Fig. 1 graphically depicts the amount of interaction between HIP1 and Huntingtin proteins with varying lengths of polyglutamine repeat.

DETAILED DESCRIPTION OF THE INVENTION

The HIP1 protein which interacts with the HD gene product was identified using the yeast two-hybrid system described in US Patent No. 5,283,173 which is incorporated herein by reference. Briefly, this system utilizes two chimeric genes or plasmids expressible in a yeast host. The yeast host is selected to contain a detectable marker gene having a binding site for the DNA binding domain of a transcriptional activator. The first chimeric gene or plasmid encodes a DNA-binding domain which recognizes the binding site of the selectable marker gene and a test protein or protein fragment. The second chimeric gene or plasmid encodes for a second test protein and a transcriptional activation domain. The two chimeric genes or plasmids are introduced into the host cell and expressed, and the cells are cultivated. Expression of the detectable marker gene only occurs when the gene product of the first chimeric gene or plasmid binds to the DNA binding domain of the detectable marker gene, and a transcriptional activation domain is brought into sufficient proximity to the DNA-binding domain, an occurrence which is facilitated by protein-protein interactions between the first and second test proteins. By selecting for cells expressing the detectable marker gene, those cells which contain chimeric genes or plasmids for interacting proteins can be identified, and the gene can be recovered and identified.

In testing for Huntington Interacting Proteins, several different plasmids were prepared containing portions of the HD gene. The first four, identified as 16pGBT9, 44pGBT9, 80pGBT9 and 128pGBT9, were GAL4 DNA binding domain-HD in-frame fusions containing nucleotides 314 to 1955 (amino acids 1-540) of the published HD cDNA sequences cloned into the vector pGBT9 (Clontech). These plasmids contain a CAG repeat region of 16, 44, 80 and 128 glutamine-encoding repeats, respectively. A clone (DMK BamHlpGBT9) was made by fusing acDNA encoding the first 544 amino acids of the myotonic dystrophy gene (a gift from R. Komeluk) in-frame with the GAL4-DNA BD of pGBT9 and was used as a negative control.

- 4 -

These plasmids have been used to identify and characterize HIP1, two additional HD-interacting proteins, HIP2 and HIP3 proteins, and can be further used for the identification of additional interacting proteins, and for tests to refine the region on the protein in which the interaction occurs. Thus, a first aspect of the invention is these four plasmids, and the use of this plasmids in identifying HD-interacting proteins. Furthermore, it will be appreciated that the GAL4 DNA-binding and activating domains are not the only domains which can be used in the yeast two-hybrid assay. Thus, in a broader sense, the invention encompasses any chimeric genes or plasmids containing nucleotides 314 to 1955 of the HD gene together with an activating or DNA-binding domain suitable for use in the yeast one, two- or three-hybrid assay for proteins critical in either binding to the HD protein or responsible for regulated expression of the HD gene.

After introducing the plasmids into Y190 yeast host cells, transforming the host cells with an adult human brain Matchmaker™ (Clontech) cDNA library coupled with a GAL4 activating domain, and selecting for the expression of two detectable marker genes to identify clones containing genes for interacting proteins, the activating domain plasmids were recovered and analyzed. As a result of this analysis, three different cDNA fragments were identified as encoding for HD-interacting proteins and designated as HIP1, HIP2 and HIP3. The sequences of HIP1 and HIP3 are given in Seq. ID Nos 1 and 3. The polypeptides which each encodes are given by Seq. ID Nos. 2 and 4. Further investigation of the HIP1 cDNA resulted in the characterization of an additional region of cDNA totaling 4795 bases and a corresponding protein, the sequences of which are given by Seq ID Nos. 5 and 6. respectively.

The cDNA molecules, particularly those encoding portions of HIP1, can be explored using oligonucleotide probes for example for amplification and sequencing. In addition, oligonucleotide probes complementary to the cDNA can be used as diagnostic probes to localize and quantify the presence of HIP1 DNA. Probes of this type with a one or two base mismatch can also be used in site-directed mutagenesis to introduce variations into the HIP1 sequence which may increase. Thus, a further aspect of the present invention is an oligonucleotide probe, preferably having a length of from 15-40 bases which specifically and selectively hybridizes with the cDNA given by Seq. ID No. 1 or 5 or a sequence complementary thereto. As used herein, the phrase "specifically and selectively hybridizes with" the

cDNA refers to primers which will hybridize with the cDNA under stringent hybridization conditions.

DNA sequencing of the HIP1 cDNA initially isolated from the yeast two-hybrid screen revealed a 1.2 kb cDNA that shows no significant degree of nucleic acid identity with any stretch of DNA using the blastn program at ncbi (blast@ncbi.nlm.nih.gov). When the entire HIP1 cDNA sequence (SEQ ID NO. 5) is translated into a polypeptide, the entire HIP1 cDNA coding (nucleotides 328-3069) is observed to be devoid of stop codons, and to produce a 914 amino acid polypeptide. A polypeptide identity search revealed an identity match over the entire length of the protein (46% conservation) with that of a hypothetical protein from *C. elegans* (ZK370.3 protein; *C. elegans* cosmid ZK370). This *C. elegans* protein shares identity with the mouse talin gene, which encodes a 217 kDa protein implicated with maintaining integrity of the cytoskeleton. It also shares identity with the SLA2/MOP2/ END4 gene from *Saccharomyces cerevisiae*, which is known to code for an essential cytoskeletal associated gene required for the accumulation and or maintenance of plasma membrane H⁺-ATPase on the cell surface. When pairwise comparisons are performed between HIP1 and the *C. elegans* ZK370.3 protein (Genpept accession number celzk370.3), it shows 26% complete identity and an overall 46% level of conservation. Comparative analysis between HIP1 and SLA2/MOP2/ END4 (EMBL accession number Z22811) demonstrate similar conservation (20% identity, 40% conservation).

HIP2 is a 2.0 kb cDNA that encodes all but the 5'-most 33 amino acids of human E2_{25k} ubiquitin conjugating enzyme. The resulting peptide has 100% identity with the previously characterized bovine E2_{25k} protein. The cDNA has 95% nucleotide identity with the bovine cDNA. Ubiquitin-conjugating enzyme is an important component in ubiquitin-mediated protein degradation pathways.

No difference in the strength of the interaction between HIP2 and HD constructs containing either 44 or 15 CAG repeats is detected using a quantitative β -galactosidase assay. The expression pattern of HIP2 (E2_{25k}) in the various parts of the brain and nervous system appears to follow the specific neuropathology observed in HD, although there does not appear to be any difference in expression levels between HD-affected and HD-non-affected individuals.

- 6 -

The third cDNA encoding an HD-interacting protein is a 537 bp cDNA coding for 187 amino acids. A search of known DNA databases did not identify the sequence homology with any known genes. However, when a protein search was performed using the blatsp server, a strong identity between HIP3 and ankyrin-related proteins was observed. The strongest identity was with the D2021.8 gene product of *C. elegans*, an uncharacterized gene, but there is also a 41% identity with AKR1, a yeast ankyrin repeat-containing protein. Furthermore, when analogous structures with charge conservation over the same amino acid stretch are considered, there is 70% protein identity. HIP3 also shares approximately 60% amino acid conservation with human brain specific ankyrins (ankyrin B and ankyrin C). Thus, it is reasonable to conclude that HIP3, like known ankyrins, is a cytoskeletal protein, and may be involved, like previously characterized ankyrins in promoting interactions between the membrane skeleton and other membrane proteins.

Further exploration of these three HD interacting proteins revealed several important facts about HIP1 that implicate it in a significantly in the pathogenesis of Huntington's Disease. First, as shown in Fig. 1, it was found that the strength of the interaction between HD protein and HIP1 is dependent on the number of CAG repeats. Second, it was found that expression of the HIP1 protein is not ubiquitous, but is limited to brain tissue. The highest amounts of expression are in the cortex, with lower levels being seen in the cerebellum, caudate and putamen.

Both HIP1 and HIP3 appear to be proteins which are involved in the maintaining the structural integrity of the cytoskeleton and various components of the cellular membrane, including microtubules and focal adhesions. Based upon this, the HD protein may be associated as part of the cytoskeletal matrix in cells where it is expressed, and our work supports the conclusion that binding of HIP1 to the HD protein is necessary for the functional incorporation of the HD protein into the cell membrane. In this circumstance, the larger polyglutamine tract in huntingtin has a decreased ability for an HIP1-HD interaction. This decreased affinity for each other disrupts the normally strong HD-HIP1-cytoskeletal anchoring association. Further, the HIP1-HD interaction may be a critical interaction at the membranes of synaptic vesicles and a decrease in the affinity of HIP1 for huntingtin may affect protein trafficking or membrane organization throughout

- 7 -

the neuron. Finally, we have demonstrated that HIP1 and HD are both found in the Triton X-100 insoluble membrane compartment of the cell, therefore, a decreased interaction between HIP1 and huntingtin may allow an abnormally subtle amount of huntingtin to be found in subcellular compartments in which it is normally found.

5 As a result of all three of these phenomenon, increased apoptosis can occur in specific neurons within the striatum. This increase in apoptosis arises from an increased susceptibility of polyglutamine-expanded huntingtin to cleavage by apopain, and because more of the expanded forms of the HD protein may be available for cleavage (and subsequent apoptosis) due to the fact they are not as tightly associated at the HD-HIP1-
10 cytoskeletal complex.

 This understanding of a biochemical basis for the pathogenesis of Huntington's Disease opens the doorway to a therapeutic method to ameliorate the pathology in patients expressing huntingtin protein with expanded polyglutamine tracts. In accordance with the method, the patient is treated to increase the amount of HIP1 or an
15 equivalent polypeptide which interacts less well with expanded Huntingtin than with Huntingtin having a CAG repeat region containing 15 to 35 repeats and facilitates the incorporation of Huntingtin into brain cell membranes.

 Increasing expression of HIP1 or an equivalent polypeptide can be accomplished using gene therapy approaches. In general, this will involve introduction of
20 DNA encoding HIP1 in an expressable vector into the brain cells. Vectors which have been shown to be suitable expression systems in mammalian cells include the herpes simplex viral based vectors: pHSV1 (Geller et al. Proc. Natl. Acad. Sci 87:8950-8954 (1990)); recombinant retroviral vectors: MFG (Jaffee et al. Cancer Res. 53:2221-2226 (1993)); Moloney-based retroviral vectors: LN, LNSX, LNCX, LXSX (Miller and
25 Rosman Biotechniques 7:980-989 (1989)); vaccinia viral vector: MVA (Sutter and Moss Proc. Natl. Acad. Sci. 89:10847-10851 (1992)); recombinant adenovirus vectors : pJM17 (Ali et al Gene Therapy 1:367-384 (1994)), (Berkner K. L. Biotechniques 6:616-624 1988); second generation adenovirus vector: DE1/DE4 adenoviral vectors (Wang and
Finer Nature Medicine 2:714-716 (1996)); and Adeno-associated viral vectors:
30 AAV/Neo (Muro-Cacho et al. J. Immunotherapy 11:231-237 (1992)).

- 8 -

Delivery of retroviral vectors to brain and nervous system tissue has been described in US Patents Nos. 4,866,042, 5,082,670 and 5,529,774, which are incorporated herein by references. These patents disclose the use of cerebral grafts or implants as one mechanism for introducing vectors bearing therapeutic gene sequences into the brain, as well as an approach in which the vectors are transmitted across the blood brain barrier.

In addition to increasing the amount of HIP1 present in brain cells of affected individuals, HD lethal phenotype may be rescued by coexpression of a HIP1 and normal sized HD protein within the same cell, specifically within neurons. The over-expression of the normal HD protein and the presence of excess HIP1 in the cell may be able to override the damaging effects of a decreased interaction between HIP1 and an expanded form of the HD protein. Therefore, a "normal state" of interaction of HD with HIP1 will rescue the cell from premature apoptotic death. Thus, a therapeutically desirable mammalian expression vector may include both a region encoding HIP1 and a region encoding normal (less than 35 repeats) HD protein.

To further illustrate the methods of making the materials which are the subject of this invention, and the testing which has established their utility, the following non-limiting experimental procedures are provided.

EXAMPLE 1

IDENTIFICATION OF INTERACTING PROTEINS

GAL4-HD cDNA constructs

An HD cDNA construct (44pGBT9), with 44 CAG repeats was generated encompassing amino acids 1 - 540 of the published HD cDNA . This cDNA fragment was fused in frame to the GAL4 DNA-binding domain (BD) of the yeast two-hybrid vector pGBT9 (Clontech). Other HD cDNA constructs, 16pGBT9, 80pGBT9 and 128pGBT9 were constructed, identical to 44pGBT9 but included only 16, 80 or 128 CAG repeats, respectively.

Another clone (DMKDBamHlpGBT9) containing the first 544 amino acids of the myotonic dystrophy gene (a gift from R. Korneluk) was fused in-frame with the

- 9 -

GAL4-DNA BD of pGBT9 and was used as a negative control. Plasmids expressing the GAL4-BDRAD7 (D. Gietz, unpublished) and SIR3 were used as a positive control for the β -galactosidase filter assay.

The clones IT15-23Q, IT15-44Q and HAP1 were generous gifts from Dr. C. Ross. These clones represent a previously isolated huntingtin interacting protein that has a higher affinity for the expanded form of the HD protein.

Yeast strains, transformations and β -galactosidase assays

The yeast strain Y190 (MATa leu2-3,112, ura3-52, trp1-901, his3- Δ 200, ade2-101, gal4 Δ gal80 Δ , URA3::GAL-lacZ, LYS2::GAL-HIS3,cyc^r) was used for all transformations and assays. Yeast transformations were performed using a modified lithium acetate transformation protocol and grown at 30 C using appropriate synthetic complete (SC) dropout media.

The β -galactosidase chromogenic filter assays were performed by transferring the yeast colonies onto Whatman filters. The yeast cells were lysed by submerging the filters in liquid nitrogen for 15-20 seconds. Filters were allowed to dry at room temperature for at least five minutes and placed onto filter paper presoaked in Z-buffer (100 mM sodium phosphate (pH7.0) 10 mM KCl, 1 mM MgSO₄) supplemented with 50 mM 2-mercaptoethanol and 0.07 mg/ml 5-bromo-4-chloro-3-indolyl β -D-galactoside (X-gal). Filters were placed at 37 C for up to 8 hours.

Yeast two-hybrid screening for huntingtin interacting protein (HIP)

cDNAs from an human adult brain Matchmaker™ cDNA library (Clontech) was transformed into the yeast strain Y190 already harboring the 44pGBT9 construct. The transformants were plated onto one hundred 150 mm x 15 mm circular culture dishes containing SC media deficient in Trp, Leu and His. The herbicide 3-amino-1,2,4-triazole (3-AT) (25mM) was utilized to limit the number of false His⁺ positives (31). The yeast transformants were placed at 30 C for 5 days and β -galactosidase filter assays were performed on all colonies found after this time, as described above, to identify β -galactosidase⁺ clones.

Primary His⁺/ β -galactosidase⁺ clones were then orderly patched onto a grid on SC

- 10 -

-Trp/-Leu/-His (25 mM 3AT) plates and assayed again for His⁺ growth and the ability to turn blue with a filter assay. Secondary positives were identified for further analysis. Proteins encoded by positive cDNAs were designated as HIPs (Huntingtin Interactive Proteins). Approximately 4.0×10^7 Trp/Leu auxotrophic transformants were screened and of 14 clones isolated 12 represented the same cDNA (HIP1), and the other 2 cDNAs, HIP2 and HIP3 were each represented only once.

The HIP cDNA plasmids were isolated by growing the His⁺/β-galactosidase⁺ colony in SC -Leu media overnight, lysing the cells with acid-washed glass beads and electroporating the bacterial strain, KC8 (leuB auxotrophic) with the yeast lysate. The KC8 ampicillin resistant colonies were replica plated onto M9 (-Leu) plates. The plasmid DNA from M9⁺ colonies was transformed into DH5-a for further manipulation.

EXAMPLE 2

CONFIRMATION OF INTERACTIONS

The HIP1-GAL4-AD cDNA activated both the lac-Z and His reporter genes in the yeast strain Y190 only when co-transformed with the GAL4-BD-HD construct, but not the negative controls (Figure 1) of the vector alone or a random fusion protein of the myotonin kinase gene. In order to assess the influence of the polyglutamine tract on the interaction between HIP1 and HD, semi-quantitative β-galactosidase assays were performed. GAL4-BD-HD fusion proteins with 16, 44, 80 and 128 glutamine repeats were assayed for their strength of interaction with the GAL4-AD-HIP1 fusion protein.

Liquid β-galactosidase assays were performed by inoculating a single yeast colony into appropriate synthetic complete (SC) dropout media and grown to OD₆₀₀ 0.6-1.5. Five millilitres of overnight culture was pelleted and washed once with 1 ml of Z-Buffer, then resuspended in 100 ml Z-Buffer supplemented with 38 mM 2-mercaptoethanol, and 0.05% SDS. Acid washed glass beads (~100 ml) were added to each sample and vortexed for four minutes, by repeatedly alternating a 30 seconds vortex, with 30 seconds on ice. Each sample was pelleted and 10 ml of lysate was added to 500 ml of lysis buffer. The samples were incubated in a 30 C waterbath for 30 seconds and then 100 ml of a 4 mg/ml o-nitrophenyl b-D galactopyranoside (ONPG) solution was added to each tube.

No difference in the β -galactosidase activity was observed between the 16 and 44 repeats or between the 80 and 128 repeats. However, a significant difference ($p < 0.05$) in activity is seen between the smaller repeats (16 and 44) and the larger repeats (80 and 128). (Figure 1)

DNA SEQUENCING, cDNA ISOLATION AND 5' RACE

Subsequently, primer walking was used to determine the remaining sequences. A human frontal cortex >4.0 kb cDNA library (a gift from S. Montal) was screened to isolate the full length HIP1 gene. Fifty nanograms of a 558 base pair Eco RI fragment from the original HIP1 cDNA was radioactively labeled with [$\alpha^{32}\text{P}$]-dCTP using nick-translation and the probe allowed to hybridized to filters containing >10⁵ pfu/ml of the cDNA library overnight at 65 C in Church buffer (see Northern blot protocol). The filters were washed at 65 C for 10 minutes with 1 X SSPE, 15 minutes at 65 C with 1 X SSPE and 0.1 % SDS, then for thirty minutes and fifteen minutes with 1 X SSPE and 0.1 % SDS. The filters were exposed to X-ray film (Kodak, XAR5) overnight at -70 C. Primary

- 12 -

positives were isolated and replated and subsequent secondary positives were hybridized and washed as for the primary screen. The resulting positive phage were converted into plasmid DNA by conventional methods (Stratagene) and the cDNA isolated and sequenced.

In order to obtain the most 5' sequence of the HIP1 gene, a Rapid Amplification of cDNA Ends (RACE) protocol was performed according to the manufacturers recommendations (BRL). First strand cDNA was synthesized using the oligo HIP1-242R (5' GCT TGA CAG TGT AGT CAT AAA GGT GGC TGC AGT CC 3'). After dCTP tailing the cDNA with terminal deoxy transferase, two rounds of 35 cycles (94 C 1 minute; 53 C 1 minute; 72 C 2 minutes) of PCR using HIP1-R2 (5' GGA CAT GTC CAG GGA GTT GAA TAC 3') and an anchor primer (5' (CUA)₄ GGC CAC GCG TCG ACT AGT ACG GGI IGG GII GGG IIG3') (BRL) were performed. The subsequent 650 base pair PCR product was cloned using the TA cloning system (Invitrogen) and sequenced using T3 and T7 primers. Sequences ID Nos. 1 and 5 show the sequence of the HIP1 cDNAs obtained.

EXAMPLE 4

DNA AND AMINO ACID ANALYSES

Overlapping DNA sequence was assembled using the program MacVector and sent via email or Netscape to the BLAST server at NIH (<http://www.ncbi.nlm.nih.gov>) to search for sequence similarities with known DNA (blastn) or protein (tblastn) sequences. Amino acid alignments were performed with the program Clustalw.

EXAMPLE 5

FISH DETECTION SYSTEM AND IMAGE ANALYSIS

The HIP1 cDNA isolated from the two-hybrid screen was mapped by fluorescent in situ hybridization (FISH) to normal human lymphocyte chromosomes counterstained with propidium iodide and DAPI. Biotinylated probe was detected with avidin-fluorescein isothiocyanate (FITC). Images of metaphase preparations were captured by a thermoelectrically cooled charge coupled camera (Photometrics). Separate images of DAPI banded chromosomes and FITC targeted chromosomes were obtained. Hybridization

signals were acquired and merged using image analysis software and pseudo colored blue (DAPI) and yellow (FITC) as described and overlaid electronically. This study showed that HIP1 maps to a single genomic locus at 7q11.2.

EXAMPLE 6

NORTHERN BLOT ANALYSIS

RNA was isolated using the single step method of homogenization in guanidinium isothiocyanate and fractionated on a 1.0% agarose gel containing 0.6 M formaldehyde. The RNA was transferred to a hybond N -membrane (Amersham) and crosslinked with ultraviolet radiation.

Hybridization of the Northern blot with b-actin as an internal control probe provided confirmation that the RNA was intact and had transferred. The 1.2 kb HIP1 cDNA was labeled using nick translation and incorporation of $\alpha^{32}\text{P}$ -dCTP. Hybridization of the original 1.2 kb HIP1 cDNA was carried out in Church buffer (0.5 M sodium phosphate buffer, pH 7.2, 2.7% sodium dodecyl sulphate, 1 mM EDTA) at 55 C overnight. Following hybridization, Northern blots were washed once for 10 minutes in 2.0 X SSPE, 0.1% SDS at room temperature and twice for 10 minutes in 0.15 X SSPE, 0.1% SDS. Autoradiography was carried out from one to three days using Hyperfilm (Amersham) film at -70 C.

Analysis of the levels of RNA levels of HIP1 by Northern blot data revealed that the 10 kilo base HIP1 message is present in all tissue assessed. However, the levels of RNA are not uniform, with brain having highest levels of expression and peripheral tissues having less message. No apparent differences in RNA expression was noted between control samples and HD affected individuals.

EXAMPLE 7

TISSUE LOCALIZATION OF HIP1

Tissue localization of HIP1 was studied using a variety of techniques as described below. Subcellular distribution of HIP-1 protein in adult human and mouse brain Biochemical fractionation studies revealed the HIP1 protein was found to be a

membrane-associated protein. No immunoreactivity was seen by Western blotting in cytosolic fractions, using the anti-HIP1-pep1 polyclonal antibody. HIP1 immunoreactivity was observed in all membrane fractions including nuclei (P1), mitochondria and synaptosomes (P2), microsomes and plasma membranes (P3). The P3 fraction contained the most HIP1 compared to other membrane fractions. HIP1 could be removed from membranes by high salt (0.5M NaCl) buffers indicating it is not an integral membrane protein, however, since low salt (0.1- 0.25M NaCl) was only able to partially remove HIP1 from membranes, its membrane association is relatively strong. The extraction of P3 membranes with the non-ionic detergent, Triton X-100 revealed HIP1 to be a Triton X-100 insoluble protein. This characteristic is shared by many cytoskeletal and cytoskeletal-associated membrane proteins including actin, which was used as a control in this study. The biochemical characteristics of HIP1 described were found to be identical in mouse and human brain and was the same for both forms of the protein (both bands of the HIP1 doublet). HIP1 co-localized with huntingtin in the P2 and P3 membrane fractions, including the high-salt membrane extractions, as well as in the Triton X-100 insoluble residue. The subcellular distribution of HIP1 was unaffected by the expression of polyglutamine-expanded huntingtin in transgenic mice and HD patient brain samples.

The localization of HIP1 protein was further investigated by immunohistochemistry in normal adult mouse brain tissue. Immunoreactivity was seen in a patchy, reticular pattern in the cytoplasm, appeared excluded from the nucleus and stained most intensely in a discontinuous pattern at the membrane. These results are consistent with the association of HIP1 with the cytoskeletal matrix and further indicate an enrichment of HIP1 at plasma membranes. Immunoreactivity occurred in all regions of the brain, including cortex, striatum, cerebellum and brainstem, but appeared most strongly in neurons and especially in cortical neurons. As described previously, huntingtin immunoreactivity was seen exclusively and uniformly in the cytosol.

The in situ hybridization studies showed HIP1 mRNA to be ubiquitously and generally expressed throughout the brain. This data is consistent with the immunohistochemical results and was identical to the distribution pattern of huntingtin mRNA in transgenic mouse brains expressing full-length human huntingtin.

Protein Preparation And Western Blotting For Expression Studies

Frozen human tissues were homogenized using a Polytron in a buffer containing 0.25M sucrose, 20mM Tris-HCl (pH 7.5), 10mM EGTA, 2mM EDTA supplemented with 10ug/ml of leupeptin, soybean trypsin inhibitor and 1mM PMSF, then centrifuged at 4,000rpm for 10' at 4 C to remove cellular debris. 100-150ug/lane of protein was separated on 8% SDS-PAGE mini-gels and then transferred to PVDF membranes. Huntingtin and HIP1 were electroblotted overnight in Towbin's transfer buffer (25 mM Tris-HCl, 0.192M glycine, pH8.3, 10% methanol) at 30V onto PVDF membranes (Immobilon-P, Millipore) as described (Towbin et al, *Proc. Nat'l Acad. Sci. (USA)* 76: 4350-4354 (1979)). Membranes were blocked for 1 hour at room temperature in 5% skim milk/ TBS (10mM Tris-HCl, 0.15M NaCl, pH7.5). Antibodies against huntingtin (pAb BKP1, 1:500), actin (mAb A-4700, Sigma, 1:500) or HIP1 (pAb HIP-pep1, 1:200) were added to blocking solution for 1 hour at room temperature. After 3 x 10 minutes washes in TBS-T (0.05% Tween-20/TBS), secondary Ab (horseradish peroxidase conjugated IgG, Biorad) was applied in blocking solution for 1 hour at room temperature. Membranes were washed and then incubated in chemiluminescent ECL solution and visualized using Hyperfilm-ECL film (Amersham).

Generation of Antibodies

The generation of huntingtin specific antibodies GHM1 and BKP1 is described elsewhere (Kalchman, et al., *J. Biol. Chem.* 271: 19385-19394 (1996)). The HIP1 peptide (VLEKDDLMDMDASQQN, a.a. 76-91 of Seq. ID No. 2) was synthesized with Cys on the N-terminus for the coupling, and coupled to Keyhole limpet hemocyanin (KLH) (Pierce) with succinimidyl 4-(N-maleimidomethyl) cyclohexane-1-carboxylate (Pierce). Female New Zealand White rabbits were injected with HIP1 peptide-KLH and Freund's adjuvant. Antibodies against the HIP1 peptide were purified from rabbit sera using affinity column with low pH elution. Affinity column was made by incubation of HIP1 peptide with activated thio-Sepharose (Pharmacia).

Western blotting of various peripheral and brain tissues were consistent with the RNA data. The HIP1 protein levels observed was not ubiquitous. The protein

- 16 -

expression is limited to brain tissue, with highest amounts seen in the cortex and lower levels seen in the cerebellum and caudate and putamen.

More regio-specific analysis of HIP1 expression in the brain revealed no differential expression pattern in affected individuals when compared to normal controls, with highest levels of expression seen in both controls and HD patients in the cortical regions.

EXAMPLE 8

CO-IMMUNOPRECIPITATION OF HIP1 WITH HUNTINGTIN

Confirmation of the HD-HIP1 interaction was performed using coimmunoprecipitation as follows. Control human brain (frontal cortex) lysate was prepared in the same manner as for subcellular localization study. Prior to immunoprecipitation, tissue lysate was centrifuged at 5000 rpm for 2 minutes at 4 C, then the supernatant was pre-cleared by the incubated with excess amount of Protein A-Sepharose for 30 minutes at 4 C, and centrifuged at the same condition. Fifty microlitres of supernatant (500 mg protein) was incubated with or without antibodies (10 ug of anti-huntingtin GHM1 (Kalchman, et al. 1996) or anti-synaptobrevin antibody) in the total 500 ul of incubation buffer (20mM Tris-Cl (pH7.5), 40mM NaCl, 1mM MgCl₂) for 1 hour at 4 C. Twenty microlitres of Protein A-Sepharose (1:1 suspension, for GHM1 and no antibody control) or Protein G-Sepharose (for anti-synaptobrevin antibody; Pharmacia) was added and incubated for 1 hour at 4 C. The beads were washed with washing buffer (incubation buffer containing 0.5 % Triton X-100) three times. The samples on the beads were separated using SDS-PAGE (7.5% acrylamide) and transferred to PVDF membrane (Immobilon-P, Millipore). The membrane was cut at about 150 kDa after transfer for Western blotting (as described above). The upper piece was probed with anti-huntingtin BKPI (1/1000) and lower piece with anti-HIP1 antibody (1/300).

The results showed that when an anti-HIP1 polyclonal antibody was immunoreacted against a blot containing the GHM1 immunoprecipitates from the brain lysate a doublet was observed at approximately 100 kDa was. When GHM1 was immunoreacted against the same immunoprecipitate the 350 kDa HD protein was also seen. The

specificity of the HD-HIP1 interaction is seen as no immunoreactive bands seen are as a result of the proteins adsorbing to the Protein-A-Sepharose (Lysate + No Antibody) or when a random, non related antibody (Lysate + anti-Synaptobrevin) is used as the immunoprecipitating antibody.

EXAMPLE 9

Subcellular fractionation of brain tissue

Cortical tissue (20-100 mg/ml) was homogenized, on ice, in a 2 ml pyrex-teflon IKA-RW15 homogenizer (Tekmar Company) in a buffer containing 0.303M sucrose, 20mM Tris-HCl pH 6.9, 1mM $MgCl_2$, 0.5mM EDTA, 1mM PMSF, 1mM leupeptin, soybean trypsin inhibitor and 1mM benzamidine (Wood et al., *Human Molec. Genet.* 5: 481-487 (1996)).

Crude membrane vesicles were isolated by two cycles of a three-step differential centrifugation protocol in a Beckman TLA 120.2 rotor at 4 C based on the methods of Wood et al (1996). The first step precipitated cellular debris and nuclei from tissue homogenates for 5 minutes at 1300 x g (P1). The 1300 x g supernatant was subsequently centrifuged for 20 minutes at 14 000 x g to isolate synaptosomes and mitochondria (P2). Finally, microsomal and plasma membrane vesicles were collected by a 35 minute centrifugation at 142 000 x g (P3). The remaining supernatant was defined as the cytosolic fraction.

High salt extraction of membranes

Aliquots of P3 membranes were twice suspended at 2mg/ ml in 0.5M NaCl, 10mM Tris-HCl, 2mM $MgCl_2$, pH7.2, containing protease inhibitors (see above). The same buffer without NaCl was used as a control. The membrane suspensions were incubated on ice for 30 minutes and then centrifuged at 142 000 x g for 30 minutes.

Extraction of cytoskeletal and cytoskeletal-associated proteins.

To extract cytoskeletal proteins, crude membrane vesicles from the P3 fraction membrane were suspended in a volume of Triton X-100 extraction buffer to give a

- 18 -

protein: detergent ratio of 5:1. The composition of the Triton X-100 extraction buffer was based on the methods of Arai et al., *J. Neuroscience* 38: 348-357 (1994) and contained 2% Triton X-100, 10mM Tris-HCl, 2mM MgCl₂, 1mM leupeptin, soybean trypsin inhibitor, PMSF and benzamidine. Membrane pellets were suspended by hand with a round-bottom teflon pestle, and placed on ice for 40 minutes. Insoluble cytoskeletal matrices were precipitated for 35 minutes at 142 000 x g in a Beckman TLA 120.2 rotor. The supernatant was defined as non-cytoskeletal-associated membrane or membrane-associated protein and was removed. The remaining pellet was extracted with Triton X-100 a second time using the same conditions. We defined the final pellet as cytoskeletal and cytoskeletal-associated protein.

Solubilization of protein and analysis by SDS-PAGE and Western Blotting

Membrane and cytoskeletal protein was solubilized in a minimum volume of 1% SDS, 3M urea, 0.1mM dithiothreitol in TBS buffer and sonicated. Protein concentration was determined using the BioRad DC Protein assay and samples were diluted at least 1 X with 5 X sample buffer (250mM Tris-HCl pH 6.8, 10% SDS, 25% glycerol, 0.02% bromophenol blue and 7% 2-mercaptoethanol) and were loaded on 7.5% SDS-PAGE gels (Bio-Rad Mini-PROTEIN II Cell system) without boiling. Western blotting was performed as described above.

Immunohistochemistry

Brain tissue was obtained from a normal C57BL/6 adult (6 months old) male mouse sacrificed with chloroform then perfusion-fixed with 4% v/v paraformaldehyde/0.01 M phosphate buffer (4% PFA). The brain tissues were removed, immersion fixed in 4% PFA for 1 day, washed in 0.01M phosphate buffered saline, pH 7.2 (PBS) for 2 days, and then equilibrated in 25% w/v sucrose PBS for 1 week. The samples were then snap-frozen in Tissue Tek molds by isopentane cooled in liquid nitrogen. After warming to -20 C, frozen blocks derived from frontal cortex, caudate/putamen, cerebellum and brainstem were cut into 14 mm sections for immunohistochemistry. Following washing in PBS, the tissue sections were blocked using 2.5% v/v normal goat serum for 1 hour at room

temperature. Primary antibodies diluted with PBS were applied to sections overnight at 4 C. Optimal dilutions for the polyclonal antibodies BKP1 and HIP1 were 1:50. Using washes of 3 x 5 minutes in PBS at room temperature, sections were sequentially incubated with biotinylated secondary antibody and then an avidin-biotin complex reagent (Vecta Stain ABC Kit, Vector) for 60 minutes each at room temperature. Color was developed using 3-3'-diaminobenzidine tetrahydrochloride and ammonium nickel sulfate.

For controls, sections were treated as described above except that HIP1 antibody aliquots were preabsorbed with an excess of HIP1 peptide as well as a peptide unrelated to HIP1 prior to incubation with the tissue sections.

In situ hybridization

In situ hybridization was performed as previously described with some modification (Suzuki et al, *BBRC* 219: 708-713 (1996)). The RNA probes were prepared using the plasmid gt149 (Lin, B., et al., *Human Molec. Genet.* 2: 1541-1545 (1994)) or a 558 subclone of HIP1. The anti-sense and sense single-stranded RNA probes were synthesized using T3 and T7 RNA polymerases and the In Vitro Transcription Kit (Clontech) with the addition of [α^{35} S]-CTP (Amersham) to the reaction mixture. Sense RNA probes were used as negative controls. For HIP1 studies normal C57BL/6 mice were used. Huntingtin probes were tested on two different transgenic mouse strains expressing full-length huntingtin, cDNA HD10366(44CAG) C57BL/6 mice and YAC HD10366(18CAG) FVB/N mice. Frozen brain sections (10um thick) were placed onto silane-coated slides under RNase-free conditions. The hybridization solution contained 40% w/v formamide, 0.02M Tris-HCl (pH 8.0), 0.005M EDTA, 0.3 M NaCl, 0.01M sodium phosphate (pH 7.0), 1x Denhardt's solution, 10% w/v dextran sulfate (pH 7.0), 0.2% w/v sarcosyl, yeast tRNA (500mg/ml) and salmon sperm DNA (200mg/ml). The radiolabelled RNA probe was added to the hybridization solution to give 1×10^6 cpm/200 ul/ section. Sections were covered with hybridization solution and incubated on formamide paper at 65 C for 18 hours. After hybridization, the slides were washed for 30 minutes sequentially with 2x SSC, 1x SSC and high stringency wash solution (50% formamide, 2x SSC and 0.1M dithiothreitol) at 65 C, followed by treatment with Rnase A

- 20 -

(1mg/ml) at 37 C for 30 minutes, then washed again and air-dried. The slides were first exposed on autoradiographic film (b-max, Amersham, UK) for 48 hours and developed for 4 minutes in Kodak D-19 followed by a 5 minute fixation in Fuji-fix. For longer exposures, the slides were dipped in autoradiographic emulsion (50% w/v in distilled water, NR-2, Konica, Japan), air-dried and exposed for 20 days at 4 C then developed as described. Sections were counterstained with methyl green or Giemsa solutions.

5

- 21 -

SEQUENCE LISTING

(1) GENERAL INFORMATION:

(i) APPLICANT: Kalchman, Michael

Goldberg, Paul

Hayden, Michael R.

(ii) TITLE OF INVENTION: Protein Which Interacts with the Huntington's Disease Gene Product, cDNA Coding Therefor, and Antibodies Therefor

(iii) NUMBER OF SEQUENCES: 8

(iv) CORRESPONDENCE ADDRESS:

(A) ADDRESSEE: Oppedahl & Larson

(B) STREET: 1992 Commerce Street Suite 309

(C) CITY: Yorktown

(D) STATE: NY

(E) COUNTRY: USA

(F) ZIP: 10598

(v) COMPUTER READABLE FORM:

(A) MEDIUM TYPE: Diskette, 3.50 inch, 1.44 Kb storage

(B) COMPUTER: IBM Compatible

(C) OPERATING SYSTEM: MS DOS 5.0

(D) SOFTWARE: WordPerfect

(vi) CURRENT APPLICATION DATA:

(A) APPLICATION NUMBER:

(B) FILING DATE:

(C) CLASSIFICATION:

(viii) ATTORNEY/AGENT INFORMATION:

(A) NAME: Larson, Marina T.

(B) REGISTRATION NUMBER: 32038

(C) REFERENCE/DOCKET NUMBER: UBC P-013

(ix) TELECOMMUNICATION INFORMATION:

(A) TELEPHONE: (914) 245-3252

(B) TELEFAX: (914) 962-4330

(2) INFORMATION FOR SEQ ID NO: 1:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 1164

(B) TYPE: nucleic acid

(C) STRANDEDNESS: single

(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: cDNA

(iii) HYPOTHETICAL: no

(iv) ANTI-SENSE: no

(vi) ORIGINAL SOURCE:

(A) ORGANISM: human

(ix) FEATURE: cDNA for Huntingtin-interacting protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO: 1:

- 22 -

ACAGCTGACA	CCCTGCAAGG	CCACCGGGAC	CGCTTCATGG	AGCAGTTTAC	50
AAAGTTGAAA	GATCTGTTCT	ACCGCTCCAG	CAACCTGCAG	TACTTCAAGC	100
GGGTCAATCA	GATCCCCCAG	CTGCCTGAGA	ACCCACCCAA	CTTCCTGCGA	150
GCCTCAGCCC	TGTCAGAACA	TATCAGCCCT	GTGGTGGTGA	TCCCTGCAGA	200
GGCCTCATCC	CCCGACAGCG	AGCCAGTCCT	AGAGAAGGAT	GACCTCATGG	250
ACATGGATGC	CTCTCAGCAG	AATTTATTTG	ACAACAAGTT	TGATGACNTC	300
TTTGGCAGTT	CATCCAGCAG	TGATCCCTTC	AATTTCAACA	GTCAAAATGG	350
TGTGAACAAG	GATGAGAAGG	ACCACTTAAT	TGAGCGACTA	TACAGAGAGA	400
TCAGTGGATT	GAAGGCACAG	CTAGAAAACA	TGAAGACTGA	GAGCCAGCGG	450
GTTGTGCTGC	AGCTGAAGGG	CCACGTCAGC	GAGCTGGAAG	CAGATCTGGC	500
CGAGCAGCAG	CACCTGCGGC	AGCAGGCGGC	CGACGACTGT	GAATTCCTGC	550
GGGCAGAACT	GGACGAGCTC	AGGNGGCAGC	GGGAGGACAC	CGAGAAGGCT	600
CAGCGGAGCC	TGTCTGAGAT	AGAAAGGAAA	GCTCAAGCCA	ATGAACAGCG	650
ATATAGCAAG	CTAAAGGAGA	AGTACAGCGA	GCTGGTTCAG	AACCACGCTG	700
ACCTGCTGCG	GAAGAATGCA	GAGGTGACCA	AACAGGTGTC	CATGGCCAGA	750
CAAGCCCAGG	TAGATTTGGA	ACGAGAGAAA	AAAGAGCTGG	AGGATTCGTT	800
GGAGCGCATC	AGTGACCAGG	GCCAGCGGAA	GACTCAAGAA	CAGCTGGAAG	850
TTCTAGAGAG	CTTGAAGCAG	GAAC TTGGCA	CAAGCCAACG	GGAGCTTCAG	900
GTTCTGCAAG	GCAGCCTGGA	AACTTCTGCC	CAGTCAGAAG	CAAAC TGGGC	950
AGCCGAGTTC	GCCGAGCTAG	AGAAGGAGCG	GGACAGCCTG	GTGAGTGCGG	1000
CAGCTCATAG	GGAGGAGGAA	TTATCTGCTC	TTCGGAAAGA	ACTGCAGGAC	1050
ACTCAGCTCA	AACTGGCCAG	CACAGAGGAA	TCTATGTGCC	AGCTTGCCAA	1100
AGACCAACGA	AAAATGCTTC	TGGTGGGGTC	CAGGAAGGCT	GCGGAGCAGG	1150
TGATACAAGA	CGCG				1164

(2) INFORMATION FOR SEQ ID NO:2:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 386

(B) TYPE: protein

(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(iii) HYPOTHETICAL: no

(vi) ORIGINAL SOURCE:

(A) ORGANISM: human

(ix) FEATURE: Huntingtin-interacting protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:2:

Thr	Ala	Asp	Thr	Leu	Gln	Gly	His	Arg	Asp	Arg	Phe	Met	Glu	Gln
1				5					10					15
Phe	Thr	Lys	Leu	Lys	Asp	Leu	Phe	Tyr	Arg	Ser	Ser	Asn	Leu	Gln
			20						25					30
Tyr	Phe	Lys	Arg	Val	Ile	Gln	Ile	Pro	Gln	Leu	Pro	Glu	Asn	Pro
			35						40					45
Pro	Asn	Phe	Leu	Arg	Ala	Ser	Ala	Leu	Ser	Glu	His	Ile	Ser	Pro
			50						55					60
Val	Val	Val	Ile	Pro	Ala	Glu	Ala	Ser	Ser	Pro	Asp	Ser	Glu	Pro
			65						70					75

- 23 -

Val	Leu	Glu	Lys	Asp	Asp	Leu	Met	Asp	Met	Asp	Ala	Ser	Gln	Gln	
				80					85					90	
Asn	Leu	Phe	Asp	Asn	Lys	Phe	Asp	Asp	Phe	Gly	Ser	Ser	Ser	Ser	
				95					100					105	
Ser	Asp	Pro	Phe	Asn	Phe	Asn	Ser	Gln	Asn	Gly	Val	Asn	Lys	Asp	
				110					115					120	
Glu	Lys	Asp	His	Leu	Ile	Glu	Arg	Leu	Tyr	Arg	Glu	Ile	Ser	Gly	
				125					130					135	
Leu	Lys	Ala	Gln	Leu	Glu	Asn	Met	Lys	Thr	Glu	Ser	Gln	Arg	Val	
				140					145					150	
Val	Leu	Gln	Leu	Lys	Gly	His	Val	Ser	Glu	Leu	Glu	Ala	Asp	Leu	
				155					160					165	
Ala	Glu	Gln	Gln	His	Leu	Arg	Gln	Gln	Ala	Ala	Asp	Asp	Cys	Glu	
				170					175					180	
Phe	Leu	Arg	Ala	Glu	Leu	Asp	Glu	Leu	Arg	Gln	Arg	Glu	Asp	Thr	
				185					190					195	
Glu	Lys	Ala	Gln	Arg	Ser	Leu	Ser	Glu	Ile	Glu	Arg	Lys	Ala	Gln	
				200					205					210	
Ala	Asn	Glu	Gln	Arg	Tyr	Ser	Lys	Leu	Lys	Glu	Lys	Tyr	Ser	Glu	
				215					220					225	
Leu	Val	Gln	Asn	His	Ala	Asp	Leu	Leu	Arg	Lys	Asn	Ala	Glu	Val	
				230					235					240	
Thr	Lys	Gln	Val	Ser	Met	Ala	Arg	Gln	Ala	Gln	Val	Asp	Leu	Glu	
				245					250					255	
Arg	Glu	Lys	Lys	Glu	Leu	Glu	Asp	Ser	Leu	Glu	Arg	Ile	Ser	Asp	
				260					265					270	
Gln	Gly	Gln	Arg	Lys	Thr	Gln	Glu	Gln	Leu	Glu	Val	Leu	Glu	Ser	
				275					280					285	
Leu	Lys	Gln	Glu	Leu	Gly	Thr	Ser	Gln	Arg	Glu	Leu	Gln	Val	Leu	
				290					295					300	
Gln	Gly	Ser	Leu	Glu	Thr	Ser	Ala	Gln	Ser	Glu	Ala	Asn	Trp	Ala	
				305					310					315	
Ala	Glu	Phe	Ala	Glu	Leu	Glu	Lys	Glu	Arg	Asp	Ser	Leu	Val	Ser	
				320					325					330	
Gly	Ala	Ala	His	Arg	Glu	Glu	Glu	Leu	Ser	Ala	Leu	Arg	Lys	Glu	
				335					340					345	

- 24 -

Leu	Gln	Asp	Thr	Gln	Leu	Lys	Leu	Ala	Ser	Thr	Glu	Glu	Ser	Met
				350					355					360
Cys	Gln	Leu	Ala	Lys	Asp	Gln	Arg	Lys	Met	Leu	Leu	Val	Gly	Ser
				365					370					375
Arg	Lys	Ala	Ala	Glu	Gln	Val	Ile	Gln	Asp	Ala				
				380					385	386				

(2) INFORMATION FOR SEQ ID NO:3:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH:

(B) TYPE: nucleic acid

(C) STRANDEDNESS: single

(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: cDNA

(iii) **HYPOTHETICAL:** no

(iv) ANTI-SENSE: no

(vi) ORIGINAL SOURCE:

(A) ORGANISM: human

(ix) **FEATURE:** cDNA for Huntingtin-interacting protein

(xi)SEQUENCE DESCRIPTION: SEQ ID NO:3:

ACCGATACCG	AAGCGGGCTG	TGTGCCCTT	CTCCACCCAG	AGGAAATCAA	50
ACCCCAAAGC	CATTATAACC	ATGGATATGG	TGAACCTCTT	GGACGGAAAA	100
CTCATATTGA	TGATTACAGC	ACATGGGACA	TAGTCAAGGC	TACACAATAT	150
GGAATATATG	AACGCTGTCG	AGAATTGGTG	GAAGCAGGTT	ATGATGTACG	200
GCAACCGGAC	AAAGAAAATG	TTACCCTCCT	CCATTGGGCT	GCCATCAATA	250
ACAGAATAGA	TTTAGTCAAA	TACTATATTT	CGAAAGGTGC	TATTGTGGAT	300
CAACTTGAG	GGGACCTGAA	TTCAACTCCA	TTGCACTGGG	ACACAAGACA	350
AGGCCATCTA	TCCATGGTTG	TGCAACTAAT	GAAATATGGT	GCAGATCCTT	400
CATTAATTGA	TGGAGAAGGA	TGTAGCTGTA	TTCATCTGGC	TGCTCAGTTC	450
GGACATACCT	CAATTGTTGC	TTATCTCATA	GCAAAAGGAC	AGGATGTG	498

(2) INFORMATION FOR SEQ ID NO:4:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 154

(B) TYPE: protein

(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(iii) **HYPOTHETICAL:** no

(vi) ORIGINAL SOURCE:

(A) ORGANISM: human

(ix) **FEATURE:** Huntingtin-interacting protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:4:

- 25 -

Thr	Asp	Thr	Glu	Ala	Gly	Cys	Val	Pro	Leu	Leu	His	Pro	Glu	Glu	1	5	10	15
Ile	Lys	Pro	Gln	Ser	His	Tyr	Asn	His	Gly	Tyr	Gly	Glu	Pro	Leu	20	25	30	
Gly	Arg	Lys	Thr	His	Ile	Asp	Asp	Tyr	Ser	Thr	Trp	Asp	Ile	Val	35	40	45	
Lys	Ala	Thr	Gln	Tyr	Gly	Ile	Tyr	Glu	Arg	Cys	Arg	Glu	Leu	Val	50	55	60	
Glu	Ala	Gly	Tyr	Asp	Val	Arg	Gln	Pro	Asp	Lys	Glu	Asn	Val	Thr	65	70	75	
Leu	Leu	His	Trp	Ala	Ala	Ile	Asn	Asn	Arg	Ile	Asp	Leu	Val	Lys	80	85	90	
Tyr	Tyr	Ile	Ser	Lys	Gly	Ala	Ile	Val	Asp	Gln	Leu	Gly	Gly	Asp	95	100	105	
Leu	Asn	Ser	Thr	Pro	Leu	His	Trp	Asp	Thr	Arg	Gln	Gly	His	Leu	110	115	120	
Ser	Met	Val	Val	Gln	Leu	Met	Lys	Tyr	Gly	Ala	Asp	Pro	Ser	Leu	125	130	135	
Ile	Asp	Gly	Glu	Gly	Cys	Ser	Cys	Ile	His	Leu	Ala	Ala	Gln	Phe	140	145	150	
Gly	His	Thr	Ser												154			

(2) INFORMATION FOR SEQ ID NO:5:

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 4846

(B) TYPE: nucleic acid

(C) STRANDEDNESS: single

(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: cDNA

(iii) HYPOTHETICAL: no

(iv) ANTI-SENSE: no

(vi) ORIGINAL SOURCE:

(A) ORGANISM: human

(ix) FEATURE: cDNA for Huntingtin-interacting protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:5:

CAGTGTACGG	TTGATCATAT	AACGCCGCGG	GCGGGGATTG	GTTTATATAT	50
CGCAAATTGA	TNTAGGGGGG	GGGGGATGGN	CAGAGATTTC	GCTTCATTAG	100
GCCATTATAA	GCAGGAAGGG	TTTCAAGGAA	AAAAACCCAG	AAAGTGCATA	150
TTGCACCCAC	CATGAGAAAG	GGGCAACAGA	CCTTNTGTTN	TGTTNTCAAC	200

- 26 -

CGCCTGCTTC	TGTTTTAGCA	ACGCAGTGTT	TTGGTGGAAG	TTGTGCCATG	250
TGTTCCACAA	ANTCTTCCGA	GATGGACACC	CGAACGTCCT	GAAGGACTTT	300
GTGAGATACA	GAAATGAATT	GAGTGACATG	AGCAGGATGT	GGGGCCACCT	350
GAGCGAGGGG	TATGGCCAGC	TGTGCAGCAT	CTACCTGAAA	CTGCTAAGAA	400
CCAAGATGGA	GTACCAACACC	AAAAATCCCA	GGTCCCAGG	CAACCTGCAG	450
ATGAGTGACC	GCCAGCTGGA	CGAGGCTGGA	GAAAGTGACG	TGAACAACCT	500
TTTCCAGTTA	ACAGTGAGGA	TGTTTGACTA	CCTGGAGTGT	GAACCTCAACC	550
TCTTCCAAAC	AGTATTCAAC	TCCCTGGACA	TGTCCCGCTC	TGTGTCCGTG	600
ACGGCAGCAG	GGCAGTGCCG	CCTCGCCCCG	CTGATCCAGG	TCATCTTGGA	650
CTGCAGCCAC	CTTTATGACT	ACACTGTCAA	GCTTCTCTTC	AAACTCCACT	700
CCTGCCTCCC	AGCTGACACC	CTGCAAGGCC	ACCGGGACCG	CTTCATGGAG	750
CAGTTTACAA	AGTTGAAAGA	TCTGTTCTAC	CGCTCCAGCA	ACCTGCAGTA	800
CTTCAAGCGG	CTCATTGAGA	TCCCCAGCT	GCCTGAGAAC	CCACCCAACCT	850
TCCTGCGAGC	CTCAGCCCTG	TCAGAACATA	TCAGCCCTGT	GGTGGTGATC	900
CCTGCAGAGG	CCTCATCCCC	CGACAGCGAG	CCAGTCCTAG	AGAAGGATGA	950
CCTCATGGAC	ATGGATGCCT	CTCAGCAGAA	TTTATTTGAC	AACAAGTTTG	1000
ATGACATCTT	TGGCAGTTCA	TTCAGCAGTG	ATCCCTTCAA	TTTCAACAGT	1050
CAAAATGGTG	TGAACAAGGA	TGAGAAGGAC	CACTTAATTG	AGCGACTATA	1100
CAGAGAGATC	AGTGGATTGA	AGGCACAGCT	AGAAAACATG	AAGACTGAGA	1150
GCCAGCGGGT	TGTGCTGCAG	CTGAAGGGCC	ACGTCAGCGA	GCTGGAAGCA	1200
GATCTGGCCG	AGCAGCAGCA	CCTGCGGCAG	CAGGCGGCCG	ACGACTGTGA	1250
ATTCCTGCGG	GCAGAACTGG	ACGAGCTCAG	GAGGCAGCGG	GAGGACACCG	1300
AGAAGGCTCA	GCGGAGCCTG	TCTGAGATAG	AAAGGAAAGC	TCAAGCCAAT	1350
GAACAGCGAT	ATAGCAAGCT	AAAGGAGAAG	TACAGCGAGC	TGGTTTCAGAA	1400
CCACGCTGAC	CTGCTGCGGA	AGAATGCAGA	GGTGACCAAA	CAGGTGTCCA	1450
TGGCCAGACA	AGCCCAGGTA	GATTTTGAAC	GAGAGAAAAA	AGAGCTGGAG	1500
GATTTCGTTG	AGCGCATCAG	TGACCAGGGC	CAGCGGAAGA	CTCAAGAACA	1550
GCTGGAAGTT	CTAGAGAGCT	TGAAGCAGGA	ACTTGGCACA	AGCCAACGGG	1600
AGCTTCAGGT	TCTGCAAGGC	AGCCTGGAAA	CTTCTGCCCA	GTCAGAAGCA	1650
AACTGGGCAG	CCGAGTTCGC	CGAGCTAGAG	AAGGAGCGGG	ACAGCCTGGT	1700
GAGTGGCGCA	GCTCATAGGG	AGGAGGAATT	ATCTGCTCTT	CGGAAAGAAC	1750
TGCAGGACAC	TCAGCTCAAA	CTGGCCAGCA	CAGAGGAATC	TATGTGCCAG	1800
CTTGCCAAAG	ACCAACGAAA	AATGCTTCTG	GTGGGGTCCA	GGAAGGCTGC	1850
GGAGCAGGTG	ATACAAGACG	CCCTGAACCA	GCTTGAAGAA	CCTCCTCTCA	1900
TCAGCTGCGC	TGGGTCTGCA	GATCAACTCC	TCTCCACGGT	CACATCCATT	1950
TCCAGCTGCA	TCGAGCAACT	GGAGAAAAGC	TGGAGCCAGT	ATCTGGCCTG	2000
CCCAGAAGAC	ATCAGTGGAC	TTCTCCATTG	CATAACCCTG	CTGGCCCACT	2050
TGACCAGCGA	CGCCATTGCT	CATGGTGCCA	CCACCTGCCT	CAGAGCCCCA	2100
CCTGAGCCTG	CCGACTCACT	GACCGAGGCC	TGTAAGCAGT	ATGGCAGGGA	2150
AACCCTCGCC	TACCTGGCCT	CCCTGGAGGA	AGAGGGAAGC	CTTGAGAATG	2200
CCGACAGCAC	AGCCATGAGG	AACTGCCTGA	GCAAGATCAA	GGCCATCGGC	2250
GAGGAGCTCC	TGCCCAGGGG	ACTGGACATC	AAGCAGGAGG	AGCTGGGGGA	2300
CCTGGTGGAC	AAGGAGATGG	CGGCCACTTC	AGCTGCTATT	GAAACTTGCA	2350
CGGCCAGAAT	AGAGGAGATG	CTCAGCAAAT	CCCGAGCAGG	AGACACAGGA	2400
GTCAAATTGG	AGGTGAATGA	AAGGATCCTT	CGTTGCTGTA	CCAGCCTCAT	2450
GCAAGCTATT	CAGGTGCTCA	TCGTGGCCTC	TAAGGACCTC	CAGAGAGAGA	2500
TTGTGGAGAG	CGGCAGGGGT	ACAGCATCCC	CTAAAGAGTT	TTATGCCAAG	2550
AACTCTCGAT	GGACAGAAGG	ACTTATCTCA	GCCTCCAAGG	CTGTGGGCTG	2600
GGGAGCCACT	GTCATGGTGG	ATGCAGCTGA	TCTGGTGGTA	CAAGGCAGAG	2650
GGAAATTTGA	GGAGCTAATG	GTGTGTTCTC	ATGAAATTGC	TGCTAGCACA	2700
GCCCAGCTTG	TGGCTGCATC	CAAGGTGAAA	GCTGATAAGG	ACAGCCCCAA	2750
CCTAGCCCAG	CTGCAGCAGG	CCTCTCGGGG	AGTGAACCAG	GCCACTGCCG	2800
GCGTTGTGGC	CTCAACCATT	TCCGGCAAAT	CACAGATCGA	AGAGACAGAC	2850
AACATGGACT	TCTCAAGCAT	GACGCTGACA	CAGATCAAAC	GCCAAGAGAT	2900

- 27 -

GGATTCTCAG	GTTAGGGTGC	TAGAGCTAGA	AAATGAATTG	CAGAAGGAGC	2950
GTCAAAAAC	GGGAGAGCTT	CGGAAAAAGC	ACTACGAGCT	TGCTGGTGTT	3000
GCTGAGGGCT	GGGAAGAAGG	AACAGAGGCA	TCTCCACCTA	CACTGCAAGA	3050
AGTGGTAACC	GAAAAAGAAT	AGAGCCAAAC	CAACACCCCA	TATGTCAGTG	3100
TAAATCCTTG	TTACCTATCT	CGTGTGTGTT	ATTTCCCCAG	CCACAGGCCA	3150
AATCCTTGGA	GTCCCAGGGG	CAGCCACACC	ACTGCCATTA	CCCAGTGCCG	3200
AGGACATGCA	TGACACTTCC	CAAAGATCCC	TCCATAGCGA	CACCCTTTCT	3250
GTTTGGACCC	ATGGTCATCT	CTGTTCTTTT	CCCCCCTCCC	TAGTTAGCAT	3300
CCAGGCTGGC	CAGTGCTGCC	CATGAGCAAG	CCTAGGTACG	AAGAGGGGTG	3350
GTGGGGGGCA	GGGCCACTCA	ACAGAGAGGA	CCAACATCCA	GTCCTGCTGA	3400
CTATTTGACC	CCCACAACAA	TGGGTATCCT	TAATAGAGGA	GCTGCTTGTT	3450
GTTTGTGAC	AGCTTGGAAA	GGGAAGATCT	TATGCCTTTT	CTTTTCTGTT	3500
TTCTTCTCAG	TCTTTTCAGT	TTCATCATTT	GCACAAACTT	GTGAGCATCA	3550
GAGGGCTGAT	GGATTCCAAA	CCAGGACACT	ACCCTGAGAT	CTGCACAGTC	3600
AGAAGGACGG	CAGGAGTGTC	CTGGCTGTGA	ATGCCAAAGC	CATTCTCCCC	3650
CTCTTTGGGC	AGTGCCATGG	ATTTCCACTG	CTTCTTATGG	TGGTTGGTTG	3700
GGTTTTTTGG	TTTTGTTTTT	TTTTTTTAAAG	TTTCACTCAC	ATAGCCAAC	3750
CTCCCAAAGG	GCACACCCCT	GGGGCTGAGT	CTCCAGGGCC	CCCCAACTGT	3800
GGTAGCTCCA	GCGATGGTGC	TGCCCAGGCC	TCTCGGTGCT	CCATCTCCGC	3850
CTCCCACTG	ACCAAGTGCT	GGCCCACCCA	GTCCATGCTC	CAGGGTCAGG	3900
CGGAGCTGCT	GAGTGACAGC	TTTCCTCAAA	AAGCAGAAGG	AGAGTGAGTG	4000
CCTTTCCCTC	CTAAAGCTGA	ATCCCGGCGG	AAAGCCTCTG	TCCGCCTTTA	4050
CAAGGGAGAA	GACAACAGAA	AGAGGGACAA	GAGGGTTCAC	ACAGCCCACT	4100
TCCCGTGACG	AGGCTCAAAA	ACTTGATCAC	ATGCTTGAAT	GGAGCTGGTG	4150
AGATCAACAA	CACTACTTCC	CTGCCGGAAT	GAAGTGTCCG	TGAATGGTCT	4200
CTGTCAAGCG	GGCCGTCTCC	CTTGGCCAG	AGACGGAGTG	TGGGAGTGAT	4250
TCCCAACTCC	TTTCTGCAGA	CGTCTGCCTT	GGCATCCTCT	TGAATAGGAA	4300
GATCGTTCCA	CTTTCTACGC	AATTGACAAA	CCCGGAAGAT	CAGATGCAAT	4350
TGCTCCCATC	AGGGAAGAAC	CCTATACTTG	GTTTGCTACC	CTTAGTATTT	4400
ATTACTAACC	TCCCTTAAGC	AGCAACAGCC	TACAAAGAGA	TGCTTGGAGC	4450
AATCAGAACT	TCAGGTGTGA	CTCTAGCAAA	GCTCATCTTT	CTGCCCGGCT	4500
ACATCAGCCT	TCAAGAATCA	GAAGAAAGCC	AAGGTGCTGG	ACTGTTACTG	4550
ACTTGGATCC	CAAAGCAAGG	AGATCATTTG	GAGCTCTTGG	GTCAGAGAAA	4600
ATGAGAAAGG	ACAGAGCCAG	CGGCTCCAAC	TCCTTTCAGC	CACATGCCCC	4650
AGGCTCTCGC	TGCCCTGTGG	ACAGGATGAG	GACAGAGGGC	ACATGAACAG	4700
CTTGCCAGGG	ATGGGCAGCC	CAACAGCACT	TTTCCTCTTC	TAGATGGACC	4750
CCAGCATTTA	AGTGACCTTC	TGATCTTGGG	AAAACAGCGT	CTTCCTTCTT	4800
TATCTATAGC	AACTCATTGG	TGGTAGCCAT	CAAGCACTTC	GGAATT	4846

(2) INFORMATION FOR SEQ ID NO:6

(i) SEQUENCE CHARACTERISTICS:

(A) LENGTH: 924

(B) TYPE: protein

(D) TOPOLOGY: linear

(ii) MOLECULE TYPE: protein

(iii) HYPOTHETICAL: no

(vi) ORIGINAL SOURCE:

(A) ORGANISM: human

(ix) FEATURE: Huntingtin-interacting protein

(xi) SEQUENCE DESCRIPTION: SEQ ID NO:6:

- 28 -

Met	Ser	Arg	Met	Trp	Gly	His	Leu	Ser	Glu	Gly	Tyr	Gly	Gln	Leu
1				5					10					15
Cys	Ser	Ile	Tyr	Leu	Lys	Leu	Leu	Arg	Thr	Lys	Met	Glu	Tyr	His
				20					25					30
Thr	Lys	Asn	Pro	Arg	Phe	Pro	Gly	Asn	Leu	Gln	Met	Ser	Asp	Arg
				35					40					45
Gln	Leu	Asp	Glu	Ala	Gly	Glu	Ser	Asp	Val	Asn	Asn	Phe	Phe	Gln
				50					55					60
Leu	Thr	Val	Glu	Met	Phe	Asp	Tyr	Leu	Glu	Cys	Glu	Leu	Asn	Leu
				65					70					75
Phe	Gln	Thr	Val	Phe	Asn	Ser	Leu	Asp	Met	Ser	Arg	Ser	Val	Ser
				80					85					90
Val	Thr	Ala	Ala	Gly	Gln	Cys	Arg	Leu	Ala	Pro	Leu	Ile	Gln	Val
				95					100					105
Ile	Leu	Asp	Cys	Ser	His	Leu	Tyr	Asp	Tyr	Thr	Val	Lys	Leu	Leu
				110					115					120
Phe	Lys	Leu	His	Ser	Cys	Leu	Pro	Ala	Asp	Thr	Leu	Gln	Gly	His
				125					130					135
Arg	Asp	Arg	Phe	Met	Glu	Gln	Phe	Thr	Lys	Leu	Lys	Asp	Leu	Phe
				140					145					150
Tyr	Arg	Ser	Ser	Asn	Leu	Gln	Tyr	Phe	Lys	Arg	Leu	Ile	Gln	Ile
				155					160					165
Pro	Gln	Leu	Pro	Glu	Asn	Pro	Pro	Asn	Phe	Leu	Arg	Ala	Ser	Ala
				170					175					180
Leu	Ser	Glu	His	Ile	Ser	Pro	Val	Val	Val	Ile	Pro	Ala	Glu	Ala
				185					190					195
Ser	Ser	Pro	Asp	Ser	Glu	Pro	Val	Leu	Glu	Lys	Asp	Asp	Leu	Met
				200					205					210
Asp	Met	Asp	Ala	Ser	Gln	Gln	Asn	Leu	Phe	Asp	Asn	Lys	Phe	Asp
				215					220					225
Asp	Ile	Phe	Gly	Ser	Ser	Phe	Ser	Ser	Asp	Pro	Phe	Asn	Phe	Asn
				230					235					240
Ser	Gln	Asn	Gly	Val	Asn	Lys	Asp	Glu	Lys	Asp	His	Leu	Ile	Glu
				245					250					255
Arg	Leu	Tyr	Arg	Glu	Ile	Ser	Gly	Leu	Lys	Ala	Gln	Leu	Glu	Asn
				260					265					270

- 29 -

Met	Lys	Thr	Glu	Ser	Gln	Arg	Val	Val	Leu	Gln	Leu	Lys	Gly	His	275	280	285
Val	Ser	Glu	Leu	Glu	Ala	Asp	Leu	Ala	Glu	Gln	Gln	His	Leu	Arg	290	295	300
Gln	Gln	Ala	Ala	Asp	Asp	Cys	Glu	Phe	Leu	Arg	Ala	Glu	Leu	Asp	305	310	315
Glu	Leu	Arg	Arg	Gln	Arg	Glu	Asp	Thr	Glu	Lys	Ala	Gln	Arg	Ser	320	325	330
Leu	Ser	Glu	Ile	Glu	Arg	Lys	Ala	Gln	Ala	Asn	Glu	Gln	Arg	Tyr	335	340	345
Ser	Lys	Leu	Lys	Glu	Lys	Tyr	Ser	Glu	Leu	Val	Gln	Asn	His	Ala	350	355	360
Asp	Leu	Leu	Arg	Lys	Asn	Ala	Glu	Val	Thr	Lys	Gln	Val	Ser	Met	365	370	375
Ala	Arg	Gln	Ala	Gln	Val	Asp	Leu	Glu	Arg	Glu	Lys	Lys	Glu	Leu	380	385	390
Glu	Asp	Ser	Leu	Glu	Arg	Ile	Ser	Asp	Gln	Gly	Gln	Arg	Lys	Thr	395	400	405
Gln	Glu	Gln	Leu	Glu	Val	Leu	Glu	Ser	Leu	Lys	Gln	Glu	Leu	Gly	410	415	420
Thr	Ser	Gln	Arg	Glu	Leu	Gln	Val	Leu	Gln	Gly	Ser	Leu	Glu	Thr	425	430	435
Ser	Ala	Gln	Ser	Glu	Ala	Asn	Trp	Ala	Ala	Glu	Phe	Ala	Glu	Leu	440	445	450
Glu	Lys	Glu	Arg	Asp	Ser	Leu	Val	Ser	Gly	Ala	Ala	His	Arg	Glu	455	460	465
Glu	Glu	Leu	Ser	Ala	Leu	Arg	Lys	Glu	Leu	Gln	Asp	Thr	Gln	Leu	470	475	480
Lys	Leu	Ala	Ser	Thr	Glu	Glu	Ser	Met	Cys	Gln	Leu	Ala	Lys	Asp	485	490	495
Gln	Arg	Lys	Met	Leu	Leu	Val	Gly	Ser	Arg	Lys	Ala	Ala	Glu	Gln	500	505	510
Val	Ile	Gln	Asp	Ala	Leu	Asn	Gln	Leu	Glu	Glu	Pro	Pro	Leu	Ile	515	520	525
Ser	Cys	Ala	Gly	Ser	Ala	Asp	His	Leu	Leu	Ser	Thr	Val	Thr	Ser	530	535	540

- 30 -

Ile	Ser	Ser	Cys	Ile	Glu	Gln	Leu	Glu	Lys	Ser	Trp	Ser	Gln	Tyr	545	550	555
Leu	Ala	Cys	Pro	Glu	Asp	Ile	Ser	Gly	Leu	Leu	His	Ser	Ile	Thr	560	565	570
Leu	Leu	Ala	His	Leu	Thr	Ser	Asp	Ala	Ile	Ala	His	Gly	Ala	Thr	575	580	585
Thr	Cys	Leu	Arg	Ala	Pro	Pro	Glu	Pro	Ala	Asp	Ser	Leu	Thr	Glu	590	595	600
Ala	Cys	Lys	Gln	Tyr	Gly	Arg	Glu	Thr	Leu	Ala	Tyr	Leu	Ala	Ser	605	610	615
Leu	Glu	Glu	Glu	Gly	Ser	Leu	Glu	Asn	Ala	Asp	Ser	Thr	Ala	Met	620	625	630
Arg	Asn	Cys	Leu	Ser	Lys	Ile	Lys	Ala	Ile	Gly	Glu	Glu	Leu	Leu	635	640	645
Pro	Arg	Gly	Leu	Asp	Ile	Lys	Gln	Glu	Glu	Leu	Gly	Asp	Leu	Val	650	655	660
Asp	Lys	Glu	Met	Ala	Ala	Thr	Ser	Ala	Ala	Ile	Glu	Thr	Cys	Thr	665	670	675
Ala	Arg	Ile	Glu	Glu	Met	Leu	Ser	Lys	Ser	Arg	Ala	Gly	Asp	Thr	680	685	690
Gly	Val	Lys	Leu	Glu	Val	Asn	Glu	Arg	Ile	Leu	Arg	Cys	Cys	Thr	695	700	705
Ser	Leu	Met	Gln	Ala	Ile	Gln	Val	Leu	Ile	Val	Ala	Ser	Lys	Asp	710	715	720
Leu	Gln	Arg	Glu	Ile	Val	Glu	Ser	Gly	Arg	Gly	Thr	Ala	Ser	Pro	725	730	735
Lys	Glu	Phe	Tyr	Ala	Lys	Asn	Ser	Arg	Trp	Thr	Glu	Gly	Leu	Ile	740	745	750
Ser	Ala	Ser	Lys	Ala	Val	Gly	Trp	Gly	Ala	Thr	Val	Met	Val	Asp	765	770	775
Ala	Ala	Asp	Leu	Val	Val	Gln	Gly	Arg	Gly	Lys	Phe	Glu	Glu	Leu	780	785	790
Met	Val	Cys	Ser	His	Glu	Ile	Ala	Ala	Ser	Thr	Ala	Gln	Leu	Val	795	800	805
Ala	Ala	Ser	Lys	Val	Lys	Ala	Asp	Lys	Asp	Ser	Pro	Asn	Leu	Ala	810	815	820

- 31 -

Gln	Leu	Gln	Gln	Ala	Ser	Arg	Gly	Val	Asn	Gln	Ala	Thr	Ala	Gly
				825					830					835
Val	Val	Ala	Ser	Thr	Ile	Ser	Gly	Lys	Ser	Gln	Ile	Glu	Glu	Thr
				840					845					850
Asp	Asn	Met	Asp	Phe	Ser	Ser	Met	Thr	Leu	Thr	Gln	Ile	Lys	Arg
				855					860					865
Gln	Glu	Met	Asp	Ser	Gln	Val	Arg	Val	Leu	Glu	Leu	Glu	Asn	Glu
				870					875					880
Leu	Gln	Lys	Glu	Arg	Gln	Lys	Leu	Gly	Glu	Leu	Arg	Lys	Lys	His
				885					890					895
Tyr	Glu	Leu	Ala	Gly	Val	Ala	Glu	Gly	Trp	Glu	Glu	Gly	Thr	Glu
				900					905					910
Ala	Ser	Pro	Pro	Thr	Leu	Gln	Glu	Val	Val	Thr	Glu	Lys	Glu	
				915					920				924	

- 32 -

CLAIMS

1. A cDNA molecule comprising the sequence given by Seq. ID No. 1.
2. A cDNA molecule comprising the sequence given by Seq. ID No. 5.
3. A polypeptide comprising the sequence given by Seq. ID. No. 2.
4. A polypeptide comprising the sequence given by Seq. ID. No. 6.
5. A chimeric gene or plasmid comprising at least nucleotides 314 to 1955 of the Huntington's Disease gene and an activating or DNA binding domain suitable for use in a yeast multi-hybrid assay.
6. The chimeric gene or plasmid according to claim 5, wherein the Huntington's Disease gene encodes a polyglutamine tract having a length of 35 or fewer residues.
7. The chimeric gene or plasmid according to claim 5, wherein the Huntington's Disease gene encodes a polyglutamine tract having a length of 36 or more residues.
8. A method for ameliorating the effects of Huntington's disease in a patient expressing Huntingtin protein with an expanded CAG repeat region, comprising the step of increasing the amount of an expressed HD-interacting polypeptide in the brain of the patient, wherein the expressed HD-interacting polypeptide interacts less well with expanded Huntingtin than with Huntingtin having a CAG repeat region containing 15 to 35 repeats and facilitates the incorporation of Huntingtin into brain cell membranes.
9. The method according to claim 8, wherein the expressed HD-interacting polypeptide comprises the sequence given by Seq. ID No. 2.

- 33 -

1 10. An antibody which binds to a polypeptide having the sequence given by
2 Seq. ID. No. 2.

1 11. The antibody of claim 10, wherein the antibody binds to amino acids
2 76-91 of the polypeptide having the sequence shown in Seq. ID No. 2.

1 12. An expression vector for expression of a gene in a mammalian host
2 comprising a region encoding an HD-interacting polypeptide, wherein the HD-interacting
3 polypeptide interacts less well with expanded Huntingtin than with Huntingtin having a CAG
4 repeat region containing 15 to 35 repeats and facilitates the incorporation of Huntingtin into
5 brain cell membranes.

1 13. An expression vector for expression of a gene in a mammalian host
2 comprising a region that is the same as or complementary to Seq. ID NO. 1.

1 14. An expression vector for expression of a gene in a mammalian host
2 comprising a region that is the same as or complementary to Seq. ID NO. 5.

1 15. The expression vector according to claims of claims 12-14, further
2 comprising a region encoding Huntingtin having a polyglutamine tract of 35 or fewer.

1 16. An oligonucleotide probe having a length of from 15-40 bases which
2 specifically and selectively hybridizes with the cDNA given by Seq. ID No. 1 or a sequence
3 complementary thereto.

1/1

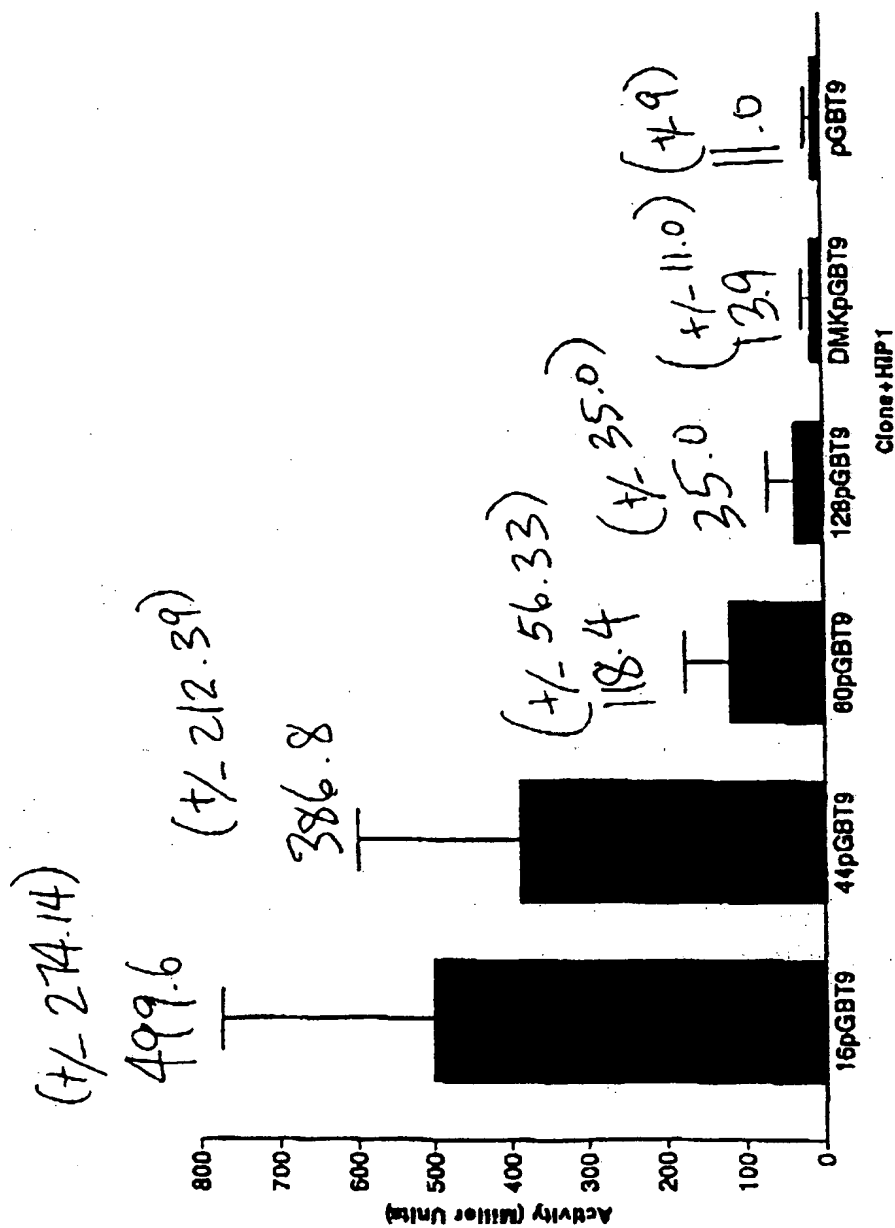


Fig. 1

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US96/18370

A. CLASSIFICATION OF SUBJECT MATTER

IPC(6) : Please See Extra Sheet.

US CL : Please See Extra Sheet.

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

U.S. : 435/320.1, 6, 69.1, 172.3; 514/44, 2; 935/62, 52, 56, 65, 34; 536/24.5, 23.1

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

STN, BIOSIS, MEDLINE, EMBASE, CAPLUS, WPIDS, APS, INPADOC

search terms: interacting protein, huntingtin, huntington, cag repeat, hip, gene therapy

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	EP 0617 125 A2 (THE GENERAL HOSPITAL CORPORATION) 28 September 1994, entire document, especially pages 4-17.	8, 10-12, 16
Y	WO 94/24279 (BERGMANN ET AL.) 27 October 1994, entire document, especially pages 13 and 28-39.	8, 12 and 16
Y	EP 0 614 977 A2 (THE GENERAL HOSPITAL CORPORATION) 14 September 1994, entire document.	1-16
Y	BIAOYANG et al. Sequence of the Murine Huntington Disease Gene: Evidence for Conservation, and polymorphism in a triplet (CCG) Repeat Alternate Splicing. Human Molecular Genetics. January 1994, Vol. 3, No. 1, pages 85-92, see entire document.	8-9, 12-15

☒ Further documents are listed in the continuation of Box C.☐ See patent family annex.

* Special categories of cited documents:	"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention.
"A" document defining the general state of the art which is not considered to be of particular relevance	"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
"E" earlier document published on or after the international filing date	"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
"L" document which may throw doubt on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	"A" document member of the same patent family
"O" document referring to an oral disclosure, use, exhibition or other means	
"P" document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search

15 JANUARY 1997

Date of mailing of the international search report

14 MAR 1997

Name and mailing address of the ISA/US
Commissioner of Patents and Trademarks
Box PCT
Washington, D.C. 20231

Facsimile No. (703) 305-3230

Authorized officer

KAREN M. HAUDA

Telephone No. (703) 308-0196

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US96/18370

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y,P	GOLDBERG et al. Cleavage of Huntingtin by apopain, a proapoptotic cysteine protease, is modulated by the polyglutamine tract. Nature Genetics. 13 August 1996, Vol. 13, No. 4, pages 442-449, see entire document.	6, 7
X,P	KALCHMAN et al. HIP-2 - A Huntingtin interacting protein: Insight into the Catabolism of the HD gene product. American Journal of Human Genetics. 02 November 1996, Vol. 59, Supplement 4, page A152, see entire document.	1-4, 12-14, and 16
Y,P		5-11 and 15

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US96/18370

Box I Observations where certain claims were found unsearchable (Continuation of item 1 of first sheet)

This international report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. ☐ Claims Nos.:
because they relate to subject matter not required to be searched by this Authority, namely:

2. ☒ Claims Nos.: 1-4, 9-11 and 13-16 (in part)
because they relate to parts of the international application that do not comply with the prescribed requirements to such an extent that no meaningful international search can be carried out, specifically:

The claims recite sequence ID numbers, but no sequence disk was submitted. Due to the length of the sequences, a search could not properly be completed on the sequence ID numbers claimed.

3. ☐ Claims Nos.:
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

Box II Observations where unity of invention is lacking (Continuation of item 2 of first sheet)

This International Searching Authority found multiple inventions in this international application, as follows:

1. ☐ As all required additional search fees were timely paid by the applicant, this international search report covers all searchable claims.
2. ☐ As all searchable claims could be searched without effort justifying an additional fee, this Authority did not invite payment of any additional fee.
3. ☐ As only some of the required additional search fees were timely paid by the applicant, this international search report covers only those claims for which fees were paid, specifically claims Nos.:

4. ☐ No required additional search fees were timely paid by the applicant. Consequently, this international search report is restricted to the invention first mentioned in the claims; it is covered by claims Nos.:

Remark on Protest

- ☐ The additional search fees were accompanied by the applicant's protest.
☐ No protest accompanied the payment of additional search fees.

INTERNATIONAL SEARCH REPORT

International application No.
PCT/US96/18370

A. CLASSIFICATION OF SUBJECT MATTER:
IPC (6):

A61K 38/00, 15/31, 15/09, 48/00; C12N 15/79, 15/63, 15/00; C07K 16/00; C07H 21/00

A. CLASSIFICATION OF SUBJECT MATTER:
US CL :

435/320.1, 6, 69.1, 172.3; 514/44, 2; 935/62, 52, 56, 65, 34; 536/24.5, 23.1